

IOWA STATE UNIVERSITY

Digital Repository

Graduate Theses and Dissertations

Iowa State University Capstones, Theses and
Dissertations

2015

Correctness results for on-line robust principal components analysis

Brian Thomas Lois

Iowa State University

Follow this and additional works at: <https://lib.dr.iastate.edu/etd>



Part of the [Electrical and Electronics Commons](#), and the [Mathematics Commons](#)

Recommended Citation

Lois, Brian Thomas, "Correctness results for on-line robust principal components analysis" (2015). *Graduate Theses and Dissertations*. 14640.

<https://lib.dr.iastate.edu/etd/14640>

This Dissertation is brought to you for free and open access by the Iowa State University Capstones, Theses and Dissertations at Iowa State University Digital Repository. It has been accepted for inclusion in Graduate Theses and Dissertations by an authorized administrator of Iowa State University Digital Repository. For more information, please contact digirep@iastate.edu.

Correctness results for on-line robust principal components analysis

by

Brian Thomas Lois

A dissertation submitted to the graduate faculty
in partial fulfillment of the requirements for the degree of
DOCTOR OF PHILOSOPHY

Co-majors: Applied Mathematics
Electrical Engineering

Program of Study Committee:
Leslie Hogben, Co-major Professor
Namrata Vaswani, Co-major Professor

Nicola Elia

Fritz Keinert

Wolfgang Kliemann

Iowa State University

Ames, Iowa

2015

DEDICATION

I dedicate this thesis to Danielle. I am very grateful for all of her love and encouragement.

TABLE OF CONTENTS

LIST OF FIGURES	vii
ACKNOWLEDGEMENTS	ix
CHAPTER 1. INTRODUCTION	1
1.1 Thesis Organization	3
References	4
CHAPTER 2. ONLINE MATRIX COMPLETION AND ONLINE ROBUST	
PCA	6
2.1 Introduction	7
2.1.1 Problem Definition	8
2.1.2 Related Work	9
2.1.3 Contributions	10
2.1.4 Notation	11
2.1.5 Organization	12
2.2 Assumptions and Main Results	12
2.2.1 Model on ℓ_t	13
2.2.2 Model on the set of missing entries or the outlier support set, \mathcal{T}_t	15
2.2.3 Denseness	17
2.2.4 Main Result for Online Matrix Completion	17
2.2.5 Main Result for Online Robust PCA	20
2.2.6 Simple Generalizations	21

2.3	Discussion	23
2.3.1	Discussion of the assumptions used	23
2.3.2	Comparison with the results for PCP and NNM	24
2.3.3	Other results for online RPCA and online MC	25
2.4	Automatic ReProCS Algorithms for Online MC and Online RPCA and Why They Work	27
2.4.1	Automatic ReProCS for Online MC (Algorithm 1)	27
2.4.2	Automatic ReProCS for online RPCA (Algorithm 2)	31
2.4.3	Key Insight for the Proof	31
2.4.4	Proof Outline	33
2.5	Most General Model on Changes in \mathcal{T}_t and a Key Lemma	34
2.5.1	Most General Model on Changes in \mathcal{T}_t	34
2.5.2	A Key Lemma that uses Model 2.5.1	35
2.6	Proof of Theorem 2.2.7 and Theorem 2.2.5	39
2.6.1	Definitions	39
2.6.2	Five Main Lemmas for Proving Theorem 2.2.7	44
2.6.3	Proof of Theorem 2.2.7	46
2.6.4	Key Lemmas for Proving of Lemmas 2.6.16, 2.6.17, and 2.6.18	48
2.6.5	Proofs of Lemmas 2.6.16, 2.6.17, and 2.6.18	49
2.7	Proofs of Lemmas 2.6.21, 2.6.22, and 2.6.23	52
2.7.1	Some definitions, remarks and facts	52
2.7.2	Preliminaries	53
2.7.3	Simple Lemmas Needed for the Proofs	54
2.7.4	Proofs of Lemma 2.6.21 and 2.6.22	55
2.7.5	Proof of Lemma 2.6.23	56
2.8	Simulation Experiments	59

2.9	Extensions	60
2.9.1	Other Models on Changes in \mathcal{T}_t	61
2.9.2	Analyze the ReProCS algorithm that also removes the deleted directions from the subspace estimate	64
2.9.3	Relax the independence assumption on ℓ_t 's	64
2.9.4	Noisy and Undersampled Online Matrix Completion or Online Robust PCA	65
2.10	Conclusions	65
2.A	Appendix A:	
	Proof that Model 2.2.3 on \mathcal{T}_t satisfies the general Model 2.5.1	66
2.B	Appendix B:	
	Proof of Lemma 2.6.14 (bound on $\zeta_{j,\text{new},k}^+$) and of Lemma 2.7.8	68
2.C	Appendix C:	
	Proof of the Compressed Sensing (CS) Lemma (Lemma 2.6.15)	70
2.D	Appendix D: Proof of Cauchy-Schwarz inequality for matrices	74
	References	74
 CHAPTER 3. RECURSIVE SPARSE RECOVERY IN CORRELATED		
	STRUCTURED NOISE	77
3.1	Introduction	77
3.1.1	Paper Organization	78
3.1.2	Problem Definition	79
3.1.3	Contribution	79
3.1.4	Notation	82
3.2	Model Assumptions, Main Result, and Discussion	83
3.2.1	Model on ℓ_t	83
3.2.2	Denseness coefficient	86
3.2.3	Model on \mathbf{x}_t	86
3.2.4	Main Result	87
3.2.5	Random Support Change	89

3.2.6	Discussion	90
3.3	Most General Support Change Model and a Key Lemma	95
3.3.1	Most General Support Change Model	95
3.3.2	Key Lemma	96
3.3.3	Unifying the signal models	97
3.4	The Automatic ReProCS Algorithm	98
3.5	Proof of Theorem 3.2.15	100
3.5.1	Definitions	100
3.5.2	Main Lemmas	102
3.5.3	Key Lemmas for Proving of Lemmas 3.5.11, 3.5.12, and 3.5.14	105
3.5.4	Proofs of Lemmas 3.5.11, 3.5.12, and 3.5.14	109
3.6	Proofs of Lemmas 3.5.20, 3.5.21, and 3.5.22	111
3.6.1	Key Lemmas Needed for the Proofs	111
3.6.2	Proofs of Lemmas 3.5.20, 3.5.21, 3.5.22	114
3.7	Proof of Lemmas 3.6.1 and 3.6.3	120
3.8	Alternative Subspace Model and Algorithm	123
3.8.1	Deletion Model	123
3.8.2	Performance Guarantee for Algorithm 4	124
3.8.3	Discussion	126
3.8.4	Cluster-PCA Algorithm (from [12])	127
3.8.5	Proof of Theorem 3.8.4	128
3.8.6	Proof of Lemma 3.8.8	132
3.9	Proofs of Lemmas 3.8.14, 3.8.15, and 3.8.16	134
3.9.1	Minor Lemmas for Proving the Main Lemmas	134
3.10	Simulations	142
3.A	Preliminaries	144
3.B	Proofs of Support Change Lemmas	151
3.C	Bounding $\zeta_{k,\text{new}}^+$ [For the purposes of review]	154
	References	155

CHAPTER 4. GENERAL CONCLUSIONS	158
References	158

LIST OF FIGURES

Figure 1.1	An fMRI sequence	2
Figure 2.1	A diagram of Model 2.2.2	13
Figure 2.3	$\varrho = 3$ and $\beta = 5$ case	16
Figure 2.4	$\varrho = 1$ and $\beta = 1$ case	16
Figure 2.5	Examples of Model 2.2.3. (a) shows a 1D video object of length s that moves by at least $s/3$ pixels once every 5 frames. (b) shows the object moving by s at every frame. (b) is an example of the best case for our result - the case with smallest ρ, β (\mathcal{T}_t 's mutually disjoint)	16
Figure 2.6	A diagram to visualize Algorithm 1 and Definition 2.6.4. The k -th projection-PCA step (at $t = \hat{t}_j + k\alpha$) computes the top left singular vectors of $(\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}') [\hat{\ell}_{\hat{t}_j + (k-1)\alpha+1}, \hat{\ell}_{\hat{t}_j + (k-1)\alpha+2}, \dots, \hat{\ell}_{\hat{t}_j + k\alpha}]$	41
Figure 2.7	Comparison of ReProCS and PCP for the RPCA problem. The top plot is the relative error $\ \ell_t - \hat{\ell}_t\ _2 / \ \ell_t\ _2$. The bottom plot shows the sparsity pattern of \mathbf{S} (black represents a non-zero entry). Results are averaged over 100 simulations and plotted every 300 time instants.	61
Figure 2.8	Model 2.9.2	62
Figure 3.2	Signal Model 3.3.1 with $\varrho = 3$	98
Figure 3.3	Disjoint Supports (Signal Model 3.3.1 with $\varrho = 1$)	98
Figure 3.4	Signal Model 3.2.12	98
Figure 3.5	Support Change Models	98
Figure 3.7	Signal Model 3.2.8	130
Figure 3.8	Algorithm 4	130

Figure 3.9	Diagrams for Signal Model 3.2.8 and Algorithm 4	130
Figure 3.10	Support of \mathbf{X} determined by Bernoulli model. The y axis is $\frac{\ \hat{\mathbf{x}}_t - \mathbf{x}_t\ _2}{\ \mathbf{x}_t\ _2}$. .	143
Figure 3.11	Support of \mathbf{X} obeys Signal Model 3.2.11. The y axis is $\frac{\ \hat{\mathbf{x}}_t - \mathbf{x}_t\ _2}{\ \mathbf{x}_t\ _2}$	144

ACKNOWLEDGEMENTS

I would like to thank Dr. Leslie Hogben for sharing her wisdom with me throughout my graduate career. Leslie has always been supportive and encouraging, helping me to achieve my goals. I would also like to thank Dr. Namrata Vaswani for her guidance and contributions to this research. Thank you also to my committee members, Dr. Nicola Elia, Dr. Fritz Keinert, and Dr. Wolfgang Kliemann, for their time and support. Finally, thank you to all of my family, friends, and fellow graduate students who have helped me complete this work.

CHAPTER 1. INTRODUCTION

This thesis considers the problem of identifying two vectors from knowledge of their sum. Obviously this is impossible in general. However, when a time series of vectors is available and some structure is assumed on both vectors, identification becomes possible. The structure assumed here is that one vector is sparse, that is, most of its entries are zero, and the other vector lies in a low-dimensional subspace, that is, if the vectors in the time series are horizontally concatenated they will form a low-rank matrix. Suppose that at each time t a vector \mathbf{m}_t is observed where

$$\mathbf{m}_t = \mathbf{x}_t + \boldsymbol{\ell}_t$$

with \mathbf{x}_t sparse and the $\boldsymbol{\ell}_t$ all belonging to some low-dimensional subspace. The goal is to recover \mathbf{x}_t and $\boldsymbol{\ell}_t$ at each time t .

Even with the assumption that the \mathbf{x}_t are sparse, and the $\boldsymbol{\ell}_t$ lie in a low-dimensional subspace, there is still an identifiability problem. Sparse vectors when horizontally concatenated could form a low rank matrix. Similarly, vectors in the low-dimensional subspace could be sparse. So it is also assumed that the \mathbf{x}_t are not low-dimensional and the $\boldsymbol{\ell}_t$ are not sparse. The property of not being sparse is called denseness, and is quantified in the sequel. Rather than directly assuming that \mathbf{x}_t are not low-rank, the results require that the supports of \mathbf{x}_t for distinct t be sufficiently different.

As a motivation for our problem, imagine video sequence with a distinct background and foreground. Suppose for example that a fixed surveillance camera records a scene as a person walks through. The background images, which do not change very much, can be modeled as belonging to a low-dimensional subspace. The foreground (person) is small compared to the size of the image, so can be modeled as a sparse vector.

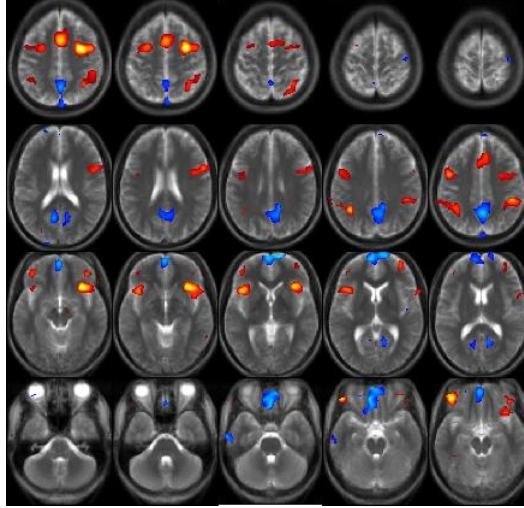


Figure 1.1: An fMRI sequence

A similar application of background/foreground separation arises in functional magnetic resonance imaging (fMRI). During an fMRI procedure, images of a patient’s brain are produced that show which areas are “active” at a given time. This can be used to map out which areas of the brain are used for different tasks. One can see from Figure 1.1 that the backgrounds from left to right remain largely similar, and the part that is lit up only covers a small amount of the overall image.

The above problem can be interpreted in two different ways. If the \mathbf{x}_t are of primary interest, the problem is one of sparse recovery. Recently, much work has been done on the problem of recovering sparse vectors. An example of early work is [1], which demonstrates the effectiveness of the ℓ_1 norm for sparse reconstruction. In [2], the author considers the following problem. Let $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b}$ where \mathbf{x} is sparse and \mathbf{b} is small noise. That is $\|\mathbf{b}\|_2 < \epsilon$ for some $\epsilon > 0$. The problem is to recover \mathbf{x} from \mathbf{y} . It is proved that under certain conditions on the matrix \mathbf{A} , the following convex program will recover \mathbf{x} with error bounded by a small constant times ϵ :

$$\min_{\tilde{\mathbf{x}}} \|\tilde{\mathbf{x}}\|_1 \text{ subject to } \|\mathbf{y} - \mathbf{A}\tilde{\mathbf{x}}\|_2 < \epsilon. \quad (1.1)$$

This procedure will be used as part of a larger algorithm to solve the original problem of separating \mathbf{x}_t and ℓ_t . The reason (1.1) alone will not work here is that ℓ_t is not necessarily small in magnitude.

Another interpretation of the original problem is as a robust principal components analysis (PCA) problem. Given a matrix of data, the goal of PCA is to compute a small number of orthogonal directions along which most of the variability in the columns lie. Because traditional procedures for PCA are sensitive to outliers in the data, much work has been done to develop algorithms for PCA that are robust with respect to outliers. In seminal work, Candès et al. [3] and Chandrasekaran et al. [4] both posed the robust PCA problem as separating a sparse matrix \mathbf{X} from a low-rank matrix \mathbf{L} from knowledge of their sum $\mathbf{M} = \mathbf{X} + \mathbf{L}$. Each proved that the convex program,

$$\min_{\tilde{\mathbf{X}}, \tilde{\mathbf{L}}} \|\tilde{\mathbf{X}}\|_{\text{sum}} + \gamma \|\tilde{\mathbf{L}}\|_* \quad \text{subject to} \quad \tilde{\mathbf{X}} + \tilde{\mathbf{L}} = \mathbf{M} \quad (1.2)$$

will exactly recover \mathbf{X} and \mathbf{L} under certain (fairly mild) conditions. In the above, γ is a scalar parameter, $\|\cdot\|_{\text{sum}}$ is the vector ℓ_1 norm of a matrix (sum of absolute values of entries), and $\|\cdot\|_*$ is the nuclear norm (sum of the singular values). The problem studied in this thesis can be viewed as a recursive or on-line version of the robust PCA problem as posed in [3] and [4]. In this thesis, the goal is to recover the columns of \mathbf{X} and \mathbf{L} as they arrive. One approach would be to re-solve (1.2) each time a new column arrives, but this is not computationally feasible when a real-time solution is desired.

This thesis analyzes versions of an algorithm called recursive projected compressed sensing (ReProCS). At a high level, ReProCS works as follows: if an accurate estimate of the subspace where the ℓ_t 's lie is available, then projecting perpendicular to this subspace estimate will nullify most of ℓ_t . Let Φ_t denote this projection. Then $\Phi_t \mathbf{m}_t = \Phi_t \mathbf{x}_t + \Phi_t \ell_t$ where $\|\Phi_t \ell_t\|_2$ is small. Finding \mathbf{x}_t is now a sparse recovery problem in small noise, so (1.1) can be used to accurately recover \mathbf{x}_t . By subtracting the estimate of \mathbf{x}_t from \mathbf{m}_t , an estimate of ℓ_t is obtained. Finally, the estimates of ℓ_t are used to maintain an accurate estimate of the subspace where the ℓ_t s lie as the subspace changes slowly over time.

1.1 Thesis Organization

The format of this thesis is journal papers in a thesis. Chapter 1 introduces the problem studied and gives some background information.

Chapter 2 contains the paper “Online Matrix Completion and Online Robust PCA,” a version of which has been submitted to *IEEE Transactions on Information Theory* and is under review. Minor modifications have been made to the submitted version for the purposes of this thesis. This paper proves a complete correctness for the ReProCS algorithm under four main assumptions: 1) an accurate estimate of the initial subspace is available; 2) the vectors ℓ_t are mutually independent over time, dense, and the subspace where they lie changes slowly; 3) the support of \mathbf{x}_t changes ‘enough’ over time; 4) the algorithm parameters are set appropriately. All of these assumptions are made precise and quantified in the paper itself. I was the primary author and researcher for this paper. The overall structure of the proof follows that in [5]. The pieces needed to obtain a complete correctness result were proved by me. Namrata Vaswani contributed both research and writing.

Chapter 3 contains the paper “Recursive Sparse Recovery in Correlated Structured Noise,” a paper that is being prepared for submission to *IEEE Transactions on Information Theory*. In ongoing work, revisions have been made, and new results have been added to the paper. The paper in this thesis contains two important improvements over the previous paper. First, the assumption that the ℓ_t s are independent over time is relaxed, and a first order autoregressive model is used for the coefficients of ℓ_t . Second, a subspace update step is added to the ReProCS algorithm. This update step allows for old directions to be deleted from the estimate of the subspace, which in turn allows for a relaxed denseness assumption on the low-dimensional vectors. I was the primary author of this paper and proved the complete correctness result. Jinchun Zhan did initial research on the correlated model, and Namrata Vaswani contributed research and writing.

Chapter 4 gives general conclusions and directions for future research.

References

- [1] S. Chen, D. Donoho, and M. Saunders, “Atomic decomposition by basis pursuit,” *SIAM Journal of Scientific Computing*, vol. 20, pp. 33–61, 1998.
- [2] E. Candes, “The restricted isometry property and its implications for compressed sensing,” *Compte Rendus de l’Academie des Sciences, Paris, Serie I*, pp. 589–592, 2008.
- [3] E. J. Candès, X. Li, Y. Ma, and J. Wright, “Robust principal component analysis?” *Journal of ACM*, vol. 58, no. 3, 2011.
- [4] V. Chandrasekaran, S. Sanghavi, P. A. Parrilo, and A. S. Willsky, “Rank-sparsity incoherence for matrix decomposition,” *SIAM Journal on Optimization*, vol. 21, 2011.
- [5] C. Qiu, N. Vaswani, B. Lois, and L. Hogben, “Recursive robust pca or recursive sparse recovery in large but structured noise,” *IEEE Trans. Info. Th.*, Aug. 2014, shorter versions in ICASSP 2013 and ISIT 2013.

CHAPTER 2. ONLINE MATRIX COMPLETION AND ONLINE ROBUST PCA

A paper submitted to *IEEE Transactions on Information Theory*

Brian Lois and Namrata Vaswani

Abstract

This work studies two interrelated problems - online robust PCA (RPCA) and online low-rank matrix completion (MC). In recent work by Candès et al., RPCA has been defined as a problem of separating a low-rank matrix (true data), $L := [\ell_1, \ell_2, \dots, \ell_t, \dots, \ell_{t_{\max}}]$ and a sparse matrix (outliers), $S := [x_1, x_2, \dots, x_t, \dots, x_{t_{\max}}]$ from their sum, $M := L + S$. Our work uses this definition of RPCA. An important application where both these problems occur is in video analytics in trying to separate sparse foregrounds (e.g., moving objects) and slowly changing backgrounds.

While there has been a large amount of recent work on both developing and analyzing batch RPCA and batch MC algorithms, the online problem is largely open. In this work, we develop a practical modification of our recently proposed algorithm to solve both the online RPCA and online MC problems. The main contribution of this work is that we obtain correctness results for the proposed algorithms under mild assumptions. The assumptions that we need are: (a) a good estimate of the initial subspace is available (easy to obtain using a short sequence of background-only frames in video surveillance); (b) the ℓ_t 's obey a 'slow subspace change' assumption; (c) the basis vectors for the subspace from which ℓ_t is generated are dense (non-sparse); (d) the support of x_t changes by at least a certain amount at least every so often; and (e) algorithm parameters are appropriately set.

2.1 Introduction

Principal Components Analysis (PCA) is a tool that is frequently used for dimension reduction. Given a matrix of data \mathbf{D} , PCA computes a small number of orthogonal directions, called principal components, that contain most of the variability of the data. For relatively noise-free data that lies close to a low-dimensional subspace, PCA is easily accomplished via singular value decomposition (SVD). The problem of PCA in the presence of outliers is referred to as robust PCA (RPCA). In recent work, Candès et al. [1] posed RPCA as a problem of separating a low-rank matrix, \mathbf{L} , and a sparse matrix, \mathbf{S} , from their sum, $\mathbf{M} := \mathbf{L} + \mathbf{S}$. They proposed a convex program called principal components' pursuit (PCP) that provided a provably correct batch solution to this problem under mild assumptions. PCP solves

$$\min_{\tilde{\mathbf{L}}, \tilde{\mathbf{S}}} \|\tilde{\mathbf{L}}\|_* + \lambda \|\tilde{\mathbf{S}}\|_{\text{sum}} \quad \text{subject to} \quad \tilde{\mathbf{L}} + \tilde{\mathbf{S}} = \mathbf{M},$$

where $\|\cdot\|_*$ is the nuclear norm (sum of singular values), $\|\cdot\|_{\text{sum}}$ is the sum of the absolute values of the entries, and λ is an appropriately chosen scalar. The same program was analyzed in parallel by Chandrasekharan et al. [2] and later by Hsu et al. [3]. Since these works, there has been a large amount of work on batch approaches for RPCA and their performance guarantees.

When RPCA needs to be solved in a recursive fashion for sequentially arriving data vectors it is referred to as online (or recursive) RPCA. Online RPCA assumes that a short sequence of outlier-free (sparse component free) data vectors is available. An example application where this problem occurs is the problem of separating a video sequence into foreground and background layers (video layering) on-the-fly [1]. Video layering is a key first step for automatic video surveillance and many other streaming video analytics tasks. In videos, the foreground usually consists of one or more moving persons or objects and hence is a sparse image. The background images usually change only gradually over time [1], e.g., moving lake waters or moving trees in a forest, and hence are well modeled as lying in a low-dimensional subspace that is fixed or slowly changing. Also, the changes are global (dense) [1]. In most video applications, it is valid to assume that an initial short sequence of background-only frames is available and this can be used to estimate the initial subspace via SVD.

Often in video applications the sparse foreground \mathbf{x}_t is actually the signal of interest, and the background ℓ_t is the noise. In this case, the problem can be interpreted as one of recursive sparse recovery in (potentially) large but structured noise. Our result allows for ℓ_t to be large in magnitude as long as it is structured. The structure we impose is that the ℓ_t 's lie in a low dimensional subspace that changes *slowly* over time.

In some other applications, instead of there being outliers, parts of a data vector may be missing entirely. When the (unknown) complete data vector is a column of a low-rank matrix, the problem of recovering it is referred to as matrix completion (MC). For example, recovering video sequences and tracking their subspace changes in the presence of easily detectable foreground occlusions. If the occluding object's intensity is known and is significantly different from that of the background, its support can be obtained by simple thresholding. The background video recovery problem then becomes an MC problem. A nuclear norm minimization (NNM) based solution for MC was introduced in [4] and studied in [5]. The convex program here is to minimize the nuclear norm of $\tilde{\mathbf{M}}$ subject to $\tilde{\mathbf{M}}$ and \mathbf{M} agreeing on all observed entries. Since then there has been a large amount of work on batch methods for MC and their correctness results.

2.1.1 Problem Definition

Consider the *online MC problem*. Let \mathcal{T}_t denote the set of missing entries at time t . We observe a vector $\mathbf{m}_t \in \mathbb{R}^n$ that satisfies

$$\mathbf{m}_t = \mathbf{I}_{\overline{\mathcal{T}}_t} \mathbf{I}_{\overline{\mathcal{T}}_t}' \ell_t \quad \text{for } t = t_{\text{train}} + 1, t_{\text{train}} + 2, \dots, t_{\text{max}}. \quad (2.1)$$

with the possibility that t_{max} can be infinity. Here ℓ_t is such that, for t large enough (quantified in Model 2.2.2), the matrix $\mathbf{L}_t := [\ell_1, \ell_2, \dots, \ell_t]$ is rank deficient. Notice that by defining \mathbf{m}_t as above, we are setting to zero the entries that are missed (see the notation section on page 11).

Consider the *online RPCA problem*. At time t we observe a vector $\mathbf{m}_t \in \mathbb{R}^n$ that satisfies

$$\mathbf{m}_t = \ell_t + \mathbf{x}_t \quad \text{for } t = t_{\text{train}} + 1, t_{\text{train}} + 2, \dots, t_{\text{max}}. \quad (2.2)$$

Here ℓ_t is as defined above and \mathbf{x}_t is the sparse (outlier) vector. We use \mathcal{T}_t to denote the support set of \mathbf{x}_t .

For both problems above, for $t = 1, 2, \dots, t_{\text{train}}$, we are given complete outlier-free measurements $\mathbf{m}_t = \ell_t$ so that it is possible to estimate the initial subspace. For the video surveillance application, this would correspond to having a short initial sequence of background only images, which can often be obtained. For $t > t_{\text{train}}$, the goal is to estimate ℓ_t (or ℓ_t and \mathbf{x}_t in case of RPCA) as soon as \mathbf{m}_t arrives and to periodically update the estimate of $\text{range}(\mathbf{L}_t)$.

In the rest of the paper, we refer to \mathcal{T}_t as the *missing/corrupted entries set*.

2.1.2 Related Work

Some other work that also studies the online MC problem (defined differently from above) includes [6, 7, 8, 9]. We discuss the connection with the idea from [6] in Section 2.4. The algorithm from [7], GROUSE, is a first order stochastic gradient method; a result for its convergence to the local minimum of the cost function it optimizes is obtained in [9]. The algorithm of [8], PETRELS, is a second order stochastic gradient method. It is shown in [8] that PETRELS converges to the stationary point of the cost function it optimizes. The advantage of PETRELS and GROUSE is that they do not need initial subspace knowledge. Another somewhat related work is [10].

Partial results have been provided for ReProCS for online RPCA in our older work [11]. In other more recent work [12] another partial result is obtained for online RPCA defined differently from above. Neither of these is a correctness result. Both require an assumption that depends on intermediate algorithm estimates. Another somewhat related work is [13] on online PCA with contaminated data. This does not model the outlier as a sparse vector but defines anything that is far from the data subspace as an outlier.

Some other works only provide an algorithm without proving any performance results, e.g., [14].

We discuss the most related works in detail in Sec 2.3.3.

2.1.3 Contributions

In this work we develop and study a practical modification of the Recursive Projected Compressive Sensing (ReProCS) algorithm introduced and studied in our earlier work [11] for online RPCA. We also develop a special case of it that solves the online MC problem. The main contribution of this work is that we obtain a *complete correctness result* for ReProCS-based algorithms for both online MC and online RPCA (or more generally, online sparse plus low-rank matrix recovery). Online algorithms are useful because they are causal (needed for applications like video surveillance) and, in most cases, are faster and need less storage compared to most batch techniques (we should mention here that there is some recent work on faster batch techniques as well, e.g., [15]). To the best of our knowledge, this work and an earlier conference version of this [16] may be among the first correctness results for online RPCA. The algorithm studied in [16] is more restrictive.

Moreover, as we will see, by exploiting temporal dependencies, such as slow subspace change, and initial subspace knowledge, our result is able to allow for a more correlated set of missing/corrupted entries than do the various results for PCP [1, 2, 3] or NNM [5] (see Sec. 2.3).

Our result uses the overall proof approach introduced in our earlier work [11] that provided a partial result for online RPCA. The most important new insight needed to get a complete result is described in Section 2.4.3. Also see Sec. 2.3.3. New proof techniques are needed for this line of work because almost all existing works only analyze batch algorithms that solve a different problem. Also, as explained in Section 2.4, the standard PCA procedure cannot be used in the subspace update step and hence results for it are not applicable.

As shown in [17], because it exploits initial subspace knowledge and slow subspace change, ReProCS has significantly improved recovery performance compared with batch RPCA algorithms - PCP [1] and [18] - as well as with the online algorithm of [14] for foreground and background extraction in many simulated and real video sequences; it is also faster than the batch methods but slower than [14].

2.1.4 Notation

We use $'$ to denote transpose. The 2-norm of a vector and the induced 2-norm of a matrix are denoted by $\|\cdot\|_2$. For a set \mathcal{T} of integers, $|\mathcal{T}|$ denotes its cardinality and $\overline{\mathcal{T}}$ denotes its complement set. We use \emptyset to denote the empty set. For a vector \mathbf{x} , $\mathbf{x}_{\mathcal{T}}$ is a smaller vector containing the entries of \mathbf{x} indexed by \mathcal{T} . Define $\mathbf{I}_{\mathcal{T}}$ to be an $n \times |\mathcal{T}|$ matrix of those columns of the identity matrix indexed by \mathcal{T} . For a matrix \mathbf{A} , define $\mathbf{A}_{\mathcal{T}} := \mathbf{A}\mathbf{I}_{\mathcal{T}}$. For matrices \mathbf{P} and \mathbf{Q} where the columns of \mathbf{Q} are a subset of the columns of \mathbf{P} , $\mathbf{P} \setminus \mathbf{Q}$ refers to the matrix of columns in \mathbf{P} and not in \mathbf{Q} .

For an $n \times n$ Hermitian matrix \mathbf{H} , $\mathbf{H} \stackrel{\text{EVD}}{=} \mathbf{U}\mathbf{\Lambda}\mathbf{U}'$ denotes an eigenvalue decomposition. That is, \mathbf{U} has orthonormal columns, and $\mathbf{\Lambda}$ is a diagonal matrix of size at least $\text{rank}(\mathbf{H}) \times \text{rank}(\mathbf{H})$. (If \mathbf{H} is rank deficient, then $\mathbf{\Lambda}$ can have any size between $\text{rank}(\mathbf{H})$ and n .) For Hermitian matrices \mathbf{A} and \mathbf{B} , the notation $\mathbf{A} \preceq \mathbf{B}$ means that $\mathbf{B} - \mathbf{A}$ is positive semi-definite. We order the eigenvalues of an Hermitian matrix in decreasing order. So $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$.

For integers a and b , we use the interval notation $[a, b]$ to mean all of the integers between a and b , inclusive, and similarly for $[a, b)$ etc.

Definition 2.1.1. For a matrix \mathbf{A} , the restricted isometry constant (RIC) $\delta_s(\mathbf{A})$ is the smallest real number δ_s such that

$$(1 - \delta_s)\|\mathbf{x}\|_2^2 \leq \|\mathbf{A}\mathbf{x}\|_2^2 \leq (1 + \delta_s)\|\mathbf{x}\|_2^2$$

for all s -sparse vectors \mathbf{x} [19]. A vector \mathbf{x} is s -sparse if it has s or fewer non-zero entries.

Definition 2.1.2. We refer to a matrix with orthonormal columns as a basis matrix. Notice that if \mathbf{P} is a basis matrix, then $\mathbf{P}'\mathbf{P} = \mathbf{I}$.

Definition 2.1.3. For basis matrices $\hat{\mathbf{P}}$ and \mathbf{P} , define $\text{dif}(\hat{\mathbf{P}}, \mathbf{P}) := \|\mathbf{I} - \hat{\mathbf{P}}\hat{\mathbf{P}}'\mathbf{P}\|_2$. This quantifies the difference between their range spaces.

If $\hat{\mathbf{P}}$ and \mathbf{P} have the same number of columns, then $\text{dif}(\hat{\mathbf{P}}, \mathbf{P}) = \text{dif}(\mathbf{P}, \hat{\mathbf{P}})$, otherwise the function is not necessarily symmetric.

2.1.5 Organization

The remainder of the paper is organized as follows. In Section 2.2 we give the model and main result for both online MC and online RPCA. Next we discuss our main results in Section 2.3. The algorithms for solving both problems are given and discussed in Section 2.4. The discussion also explains why the proof of our main result should go through. Section 2.4.3 within this section describes the key insight needed by the proof and Section 2.4.4 gives the proof outline. The most general form of our model on the missing entries set, \mathcal{T}_t , is described in Section 2.5. A key new lemma for proving our main results is also given in this section. The proof of our main results can be found in Section 2.6. Proofs of three long lemmas needed for proving the lemmas leading to the main theorem are postponed until Section 2.7. Section 2.8 shows numerical experiments backing up our claims. We discuss some extensions in Section 2.9 and give conclusions in Section 2.10

2.2 Assumptions and Main Results

Before we give our model on ℓ_t , we need the following definition.

Definition 2.2.1. Recall that $\mathbf{m}_t = \ell_t$ for $t = 1, \dots, t_{\text{train}}$ is the training data. Let $\hat{\lambda}_{\text{train}}^-$ be the minimum non-zero eigenvalue of $\frac{1}{t_{\text{train}}} \sum_{t=1}^{t_{\text{train}}} \mathbf{m}_t \mathbf{m}_t'$. That is

$$\hat{\lambda}_{\text{train}}^- := \min_{\lambda_i > 0} \lambda_i \left(\frac{1}{t_{\text{train}}} \sum_{t=1}^{t_{\text{train}}} \mathbf{m}_t \mathbf{m}_t' \right)$$

Define $\hat{\mathbf{P}}_{t_{\text{train}}}$ to be the matrix containing the eigenvectors of $\frac{1}{t_{\text{train}}} \sum_{t=1}^{t_{\text{train}}} \mathbf{m}_t \mathbf{m}_t'$, with eigenvalues larger than or equal to $\hat{\lambda}_{\text{train}}^-$, as its columns.

We will use $\hat{\lambda}_{\text{train}}^-$ in our algorithms to set the eigenvalue threshold to both detect subspace change and estimate the number of newly added directions. We also use $\hat{\lambda}_{\text{train}}^-$ to state the slow subspace change assumption below. We will use $\hat{\mathbf{P}}_{t_{\text{train}}}$ as the initial subspace knowledge in the algorithms.

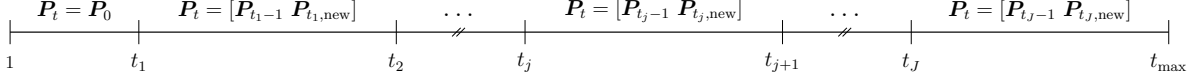


Figure 2.1: A diagram of Model 2.2.2

2.2.1 Model on ℓ_t

We assume that ℓ_t is a vector from a slowly changing low-dimensional subspace that changes in such a way that the matrix $\mathbf{L}_t := [\ell_1, \ell_2, \dots, \ell_t]$ is low-rank for t large enough. This can be modeled in various ways. One possible model is given below. It assumes that ℓ_t 's are zero mean, bounded and mutually independent random variables with a covariance matrix that is low-rank at each time and that changes “slowly” in the following fashion: (a) its column subspace remains constant for a long enough time and then changes; (b) when it changes, the number of newly added directions is small and the eigenvalues along the newly added directions are small for d frames after the change.

Model 2.2.2 (Model on ℓ_t). *Assume that the ℓ_t are zero mean and bounded random vectors in \mathbb{R}^n that are mutually independent over time. Also assume that their covariance matrix Σ_t has an eigenvalue decomposition*

$$\mathbb{E}[\ell_t \ell_t'] = \Sigma_t \stackrel{\text{EVD}}{=} \mathbf{P}_t \mathbf{\Lambda}_t \mathbf{P}_t'$$

where \mathbf{P}_t changes as

$$\mathbf{P}_t = \begin{cases} [\mathbf{P}_{t-1} \ \mathbf{P}_{t,\text{new}}] & \text{if } t = t_1 \text{ or } t_2 \text{ or } \dots \text{ or } t_J \\ \mathbf{P}_{t-1} & \text{otherwise.} \end{cases} \quad (2.3)$$

and $\mathbf{\Lambda}_t$ changes as follows. For $t \in [t_j, t_{j+1})$, define $\mathbf{\Lambda}_{t,\text{new}} := \mathbf{P}_{t_j,\text{new}}' \Sigma_t \mathbf{P}_{t_j,\text{new}}$ and assume that

$$(\mathbf{\Lambda}_{t,\text{new}})_{i,i} = (v_i)^{t-t_j} q_i \hat{\lambda}_{\text{train}}^- \text{ for } i = 1, \dots, r_{j,\text{new}} \quad (2.4)$$

where $q_i \geq 1$ and $v_i > 1$ but not too large¹. We assume that (a) $t_{j+1} - t_j \geq d$ for a $d \geq (K+2)\alpha$;

¹Our result would still hold if the v_i were different for each change time (i.e. $v_{j,i}$). We let them be the same to reduce notation.

and (b) for all i , $q_i(v_i)^d \leq 3$. Here K and α are algorithm parameters that are set in Theorem 2.2.7.

Other minor assumptions are as follows. (i) Define $t_0 := 1$ and assume that $t_{\text{train}} \in [t_0, t_1]$. (ii) For $j = 0, 1, 2, \dots, J$, define

$$r_j := \text{rank}(\mathbf{P}_{t_j}) \quad \text{and} \quad r_{j,\text{new}} := \text{rank}(\mathbf{P}_{t_{j,\text{new}}}).$$

and assume that $r_J < \min(n, t_{j+1} - t_j)$. This ensures that, for all $t > r_J$, the matrix \mathbf{L}_t is low-rank. (iii) Define

$$\lambda^+ := \sup_t \lambda_{\max}(\mathbf{\Lambda}_t)$$

as the maximum eigenvalue at any time and assume that $\lambda^+ < \infty$.

Observe from the above that \mathbf{P}_t is a basis matrix and $\mathbf{\Lambda}_t$ is diagonal. We refer to the t_j 's as the subspace change times.

A visual depiction of the above model can be found in Figure 2.1.

Define the largest and smallest eigenvalues along the new directions for the first d frames after a subspace change as

$$\lambda_{\text{new}}^+ := \max_j \max_{t \in [t_j, t_j + d]} \lambda_{\max}(\mathbf{\Lambda}_{t,\text{new}}) \quad \text{and} \quad \lambda_{\text{new}}^- := \min_j \min_{t \in [t_j, t_j + d]} \lambda_{\min}(\mathbf{\Lambda}_{t,\text{new}})$$

The slow change model on $\mathbf{\Lambda}_{t,\text{new}}$ is one way to ensure that

$$\hat{\lambda}_{\text{train}}^- \leq \lambda_{\text{new}}^- \leq \lambda_{\text{new}}^+ \leq 3\hat{\lambda}_{\text{train}}^- \quad (2.5)$$

i.e. the maximum variance of the projection of ℓ_t along the new directions is small enough for the first d frames after a change. Also the minimum variance is larger than a constant greater than zero (and hence detectable). The proof of our main result only relies on (2.5) and does not use the actual slow increase model in any other way. The above inequality along with $t_{j+1} - t_j \geq d \geq (K + 2)\alpha$ quantifies “slow subspace change”.

Notice that the above model does not put any assumption on the eigenvalues along the existing directions. In particular, they do not need to be greater than zero and hence the model automatically allows existing directions (columns of $\mathbf{P}_{t_{j-1}}$ for $t \in [t_j, t_{j+1})$) to drop out of the current subspace. It could be the case that for some time period, $(\mathbf{\Lambda}_t)_{i,i} = 0$ (for an

i corresponding to a column of \mathbf{P}_{t_j-1}), so that the i^{th} column of \mathbf{P}_{t_j-1} is not contributing anything to ℓ_t at that time. For the same index i , $(\mathbf{\Lambda}_t)_{i,i}$ could also later increase again to a nonzero value. Therefore $r_0 + \sum_{i=1}^j r_{i,\text{new}}$ is only a bound on the rank of $\mathbf{\Sigma}_t$ for $t \in [t_j, t_{j+1})$, and not necessarily the rank itself. A more explicit model for deletion of directions is to let \mathbf{P}_t change as

$$\mathbf{P}_t = \begin{cases} [(\mathbf{P}_{t-1} \setminus \mathbf{P}_{t,\text{del}}) \ \mathbf{P}_{t,\text{new}}] & \text{if } t = t_1 \text{ or } t_2 \text{ or } \dots \ t_J \\ \mathbf{P}_{t-1} & \text{otherwise.} \end{cases} \quad (2.6)$$

where $\mathbf{P}_{t,\text{del}}$ contains the columns of \mathbf{P}_{t-1} for which the variance is zero. If we add the assumption that $[\mathbf{P}_{t_1-1} \ \mathbf{P}_{t_1,\text{new}} \ \mathbf{P}_{t_2,\text{new}} \ \dots \ \mathbf{P}_{t_J,\text{new}}]$ be a basis matrix (i.e. deleted directions cannot be part of a later $\mathbf{P}_{t_j,\text{new}}$), then this is a special case of Model 2.2.2 above. We say special case because this only allows deletions at times t_j , whereas Model 2.2.2 allows deletion of old directions at any time.

For $t \in [t_j, t_{j+1})$, let $\mathbf{P}_{t,*} := \mathbf{P}_{t_j-1}$ and $\mathbf{\Lambda}_{t,*} := \mathbf{P}_{t,*}' \mathbf{\Sigma}_t \mathbf{P}_{t,*}$. Observe that Model 2.2.2 does not have any constraint on $\mathbf{\Lambda}_{t,*}$. Thus if we assume that its entries are such that their changes from t to $t+1$ are smaller than or equal to $\|\mathbf{\Lambda}_{t,\text{new}} - \mathbf{\Lambda}_{t+1,\text{new}}\|_2$, then clearly, $\frac{\|\mathbf{\Sigma}_{t+1} - \mathbf{\Sigma}_t\|_2}{\|\mathbf{\Sigma}_t\|_2} \leq (3^{1/d} - 1)$ for all $t \in [t_j, t_j + d]$ and all j ². Since d is large, the upper bound is a small quantity, i.e. the covariance matrix changes slowly. For later time instants, we do not have any requirement (and so in particular $\mathbf{\Sigma}_t$ could still change slowly). Hence the above model includes “slow changing” and low-rank $\mathbf{\Sigma}_t$ as a special case.

2.2.2 Model on the set of missing entries or the outlier support set, \mathcal{T}_t

Our result requires that the set of missing entries (or the outlier support sets), \mathcal{T}_t , have *some* changes over time. We give one simple model for it below. One example that satisfies this model is a video application consisting of a foreground with one object of length s or less that can remain static for at most β frames at a time. When it moves, it moves *downwards* (or upwards, but always in one direction) by at least s/ρ pixels, and at most s/ρ_2 pixels. Once

²This follows because $\|\mathbf{\Sigma}_t\|_2 \geq \|\mathbf{\Lambda}_{t,\text{new}}\|_2 = \max_i (v_i)^{t-t_j} q_i \hat{\lambda}_{\text{train}}^-$ and $\|\mathbf{\Sigma}_{t+1} - \mathbf{\Sigma}_t\|_2 \leq \|\mathbf{\Lambda}_{t+1,\text{new}} - \mathbf{\Lambda}_{t,\text{new}}\|_2 \leq \max_i (v_i)^{t-t_j} q_i \hat{\lambda}_{\text{train}}^- (v_i - 1) \leq \max_i (v_i)^{t-t_j} q_i \hat{\lambda}_{\text{train}}^- \max_i (v_i - 1)$. Thus the ratio is bounded by $\max_i (v_i - 1) \leq (3/q_i)^{1/d} - 1 < (3^{1/d} - 1)$ since $q_i \geq 1$.

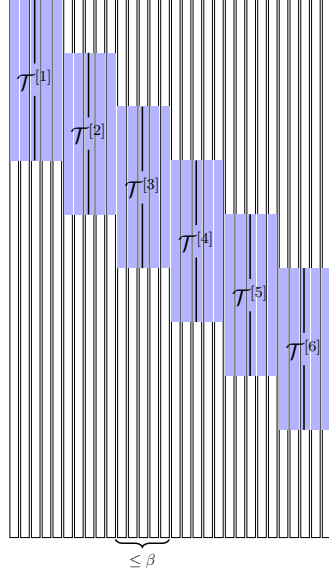
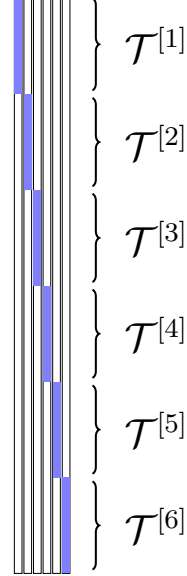
Figure 2.3 $\varrho = 3$ and $\beta = 5$ caseFigure 2.4 $\varrho = 1$ and $\beta = 1$ case

Figure 2.5: Examples of Model 2.2.3. (a) shows a 1D video object of length s that moves by at least $s/3$ pixels once every 5 frames. (b) shows the object moving by s at every frame. (b) is an example of the best case for our result - the case with smallest ρ, β (\mathcal{T}_t 's mutually disjoint)

it reaches the bottom of the scene, it disappears. The maximum motion is such that, if the object were to move at each frame, it still does not go from the top to the bottom of the scene in a time interval of length α , i.e. $\frac{s}{\rho_2}\alpha \leq n$. Anytime after it has disappeared another object could appear. We show this example in Fig. 2.5.

Model 2.2.3 (model on \mathcal{T}_t). Let t^k , with $t^k < t^{k+1}$, denote the times at which \mathcal{T}_t changes and let $\mathcal{T}^{[k]}$ denote the distinct sets. For an integer α (we set its value in Theorem 2.2.7), assume the following.

1. Assume that $\mathcal{T}_t = \mathcal{T}^{[k]}$ for all times $t \in [t^k, t^{k+1})$ with $(t^{k+1} - t^k) \leq \beta$ and $|\mathcal{T}^{[k]}| \leq s$.
2. Let ρ be a positive integer so that for any k ,

$$\mathcal{T}^{[k]} \cap \mathcal{T}^{[k+\rho]} = \emptyset;$$

assume that

$$\rho^2 \beta \leq 0.01 \alpha.$$

3. For any k ,

$$\sum_{i=k+1}^{k+\alpha} |\mathcal{T}^{[i]} \setminus \mathcal{T}^{[i+1]}| \leq n$$

and for any $k < i \leq k + \alpha$,

$$(\mathcal{T}^{[k]} \setminus \mathcal{T}^{[k+1]}) \cap (\mathcal{T}^{[i]} \setminus \mathcal{T}^{[i+1]}) = \emptyset.$$

(One way to ensure $\sum_{i=k+1}^{k+\alpha} |\mathcal{T}^{[i]} \setminus \mathcal{T}^{[i+1]}| \leq n$ is to require that for all i , $|\mathcal{T}^{[i]} \setminus \mathcal{T}^{[i+1]}| \leq \frac{s}{\rho_2}$

with $\frac{s}{\rho_2} \alpha \leq n$.)

In this model, k takes values $1, 2, \dots$; the largest value it can take is t_{\max} (this will happen if \mathcal{T}_t changes at every time).

Clearly the video moving object example satisfies the above model as long as $\rho^2 \beta \leq 0.01 \alpha$.

³ This becomes clearer from Fig. 2.5.

2.2.3 Denseness

In order to recover the ℓ_t 's from missing data or to separate them from the sparse outliers, the basis vectors for the subspace from which they are generated cannot be sparse. We quantify this using the incoherence condition from [1]. Let μ be the smallest real number such that

$$\max_i \|\mathbf{P}_{t_0}' \mathbf{I}_i\|_2^2 \leq \frac{\mu r_0}{n} \quad \text{and} \quad \max_i \|\mathbf{P}_{t_j, \text{new}}' \mathbf{I}_i\|_2^2 \leq \frac{\mu r_{j, \text{new}}}{n} \quad \text{for all } j \quad (2.7)$$

Recall from the notation section that \mathbf{I}_i is the i^{th} column of the identity matrix (or i^{th} standard basis vector). We bound μr_0 and $\mu r_{j, \text{new}}$ in the theorem.

2.2.4 Main Result for Online Matrix Completion

Definition 2.2.4. Recall that $r_{j, \text{new}} := \text{rank}(\mathbf{P}_{t_j, \text{new}})$ and $r_j := \text{rank}(\mathbf{P}_{t_j})$. Define $r_{\text{new}} := \max_j r_{j, \text{new}}$, and $r = r_0 + J r_{\text{new}}$.

³Let \mathcal{T}_t be the support set of the object (set of pixels containing the object). The first condition holds since there is at most one object of size s or less and the object cannot remain static for more than β frames. Since it moves in one direction by at least s/ρ each time it moves, this means that definitely after it moves ρ times, the supports will be disjoint (second condition). The third condition holds because it moves in one direction and by at most s/ρ_2 with $\frac{s}{\rho_2} \alpha \leq n$ (so even if it were to move at each t , i.e. if $t_{k+1} = t_k + 1$ for all k , the third condition will hold). Also see Fig. 2.5.

Also define $\mathbf{a}_t := \mathbf{P}_t' \boldsymbol{\ell}_t$, and for $t \in [t_j, t_{j+1})$, $\mathbf{a}_{t,\text{new}} := \mathbf{P}_{t_j,\text{new}}' \boldsymbol{\ell}_t$. Let

$$\gamma := \max_t \|\mathbf{a}_t\|_\infty \quad \text{and} \quad \gamma_{\text{new}} := \max_j \max_{t \in [t_j, t_j+d]} \|\mathbf{a}_{t,\text{new}}\|_\infty$$

Notice that $\text{rank}(\mathbf{L}) = \text{rank}(\mathbf{P}_{t_{\max}}) \leq r$. Also, $\|\mathbf{a}_t\|_2 \leq \sqrt{r}\gamma$ and for $t \in [t_j, t_j + d]$, $\|\mathbf{a}_{t,\text{new}}\|_2 \leq \sqrt{r_{\text{new}}}\gamma_{\text{new}}$.

The following theorem gives a correctness result for Algorithm 1 given and explained in Section 2.4. The algorithm has two parameters - α and K . The parameter α is the number of consecutive time instants that are used to obtain an estimate of the new subspace, and K is the total number of times the new subspace is estimated before we get an accurate enough estimate of it. The algorithm uses $\hat{\lambda}_{\text{train}}^-$ and $\hat{\mathbf{P}}_{t_{\text{train}}}$ defined in Definition 2.2.1 and \mathbf{m}_t as inputs.

Theorem 2.2.5. *Consider Algorithm 1. Assume that \mathbf{m}_t satisfies (2.1). Pick a ζ that satisfies*

$$\zeta \leq \min \left\{ \frac{10^{-4}}{r^2}, \frac{0.03\hat{\lambda}_{\text{train}}^-}{r^2\lambda^+}, \frac{1}{r^3\gamma^2}, \frac{\hat{\lambda}_{\text{train}}^-}{r^3\gamma^2} \right\}.$$

Suppose that the following hold.

1. $\text{dif}(\hat{\mathbf{P}}_{t_{\text{train}}}, \mathbf{P}_{t_{\text{train}}}) \leq r_0\zeta$ (notice from Model 2.2.2 that $\mathbf{P}_{t_{\text{train}}} = \mathbf{P}_{t_0} = \mathbf{P}_1$);

2. The algorithm parameters are set as:

$$K = \left\lceil \frac{\log(0.16r_{\text{new}}\zeta)}{\log(0.83)} \right\rceil; \text{ and } \alpha = C(\log(6(K+1)J) + 11\log(n)) \text{ for a constant}$$

$$C \geq C_{\text{add}} := 32 \cdot 100^2 \frac{\max\{16, 1.2(\sqrt{\zeta} + \sqrt{r_{\text{new}}}\gamma_{\text{new}})^4\}}{(r_{\text{new}}\zeta\hat{\lambda}_{\text{train}}^-)^2}; \quad (2.8)$$

3. (Subspace change) Model 2.2.2 on $\boldsymbol{\ell}_t$ holds;

4. (Changes in the missing/corrupted sets \mathcal{T}_t) Model 2.2.3 on \mathcal{T}_t holds or its generalization, Model 2.5.1 (given in Section 2.5), holds;

5. (Denseness and bound on s , r_0 , r_{new}) the bounds in (2.7) hold with $2s(r_0 + Jr_{\text{new}})\mu \leq 0.09n$ and $2sr_{\text{new}}\mu \leq 0.0004n$;

Then, with probability at least $1 - n^{-10}$, at all times t ,

$$1. \|\hat{\boldsymbol{\ell}}_t - \boldsymbol{\ell}_t\|_2 \leq 1.2(\sqrt{\zeta} + \sqrt{r_{\text{new}}}\gamma_{\text{new}})$$

2. the subspace error $\text{SE}_t := \|(\mathbf{I} - \hat{\mathbf{P}}_t \hat{\mathbf{P}}_t') \mathbf{P}_t\|_2$ is bounded above by $10^{-2} \sqrt{\zeta}$ for $t \in [t_j + d, t_{j+1})$.

Proof. The proof is given in Sections 2.6 and 2.7. As shown in Lemma 2.5.2, Model 2.2.3 is a special case of Model 2.5.1 (Model 2.5.1 is more general) on \mathcal{T}_t . Hence we prove the result only using Model 2.5.1. \square

Theorem 2.2.5 says that if an accurate estimate of the initial subspace is available; the two algorithm parameters are set appropriately; the ℓ_t 's are mutually independent over time and the low-dimensional subspace from which ℓ_t is generated changes “slowly” enough, i.e. (a) the delay between change times is large enough ($d \geq (K + 2)\alpha$) and (b) the eigenvalues along the newly added directions are small enough for d frames after a subspace change (so that (3b) holds); the set of missing entries at time t , \mathcal{T}_t , has enough changes; and the basis vectors that span the low-dimensional subspaces are dense enough; then, with high probability (w.h.p.), the error in estimating ℓ_t will be small at all times t . Also, the error in estimating the low-dimensional subspace will be initially large when new directions are added, but will decay to a small constant times $\sqrt{\zeta}$ within a finite delay.

Consider the accurate initial subspace assumption. If the training data truly satisfies $\mathbf{m}_t = \ell_t$ (without any noise or modeling error) and if we have at least r_0 linearly independent ℓ_t 's (if ℓ_t 's are continuous random vectors, this corresponds to needing $t_{\text{train}} \geq r_0$ almost surely), then the estimate of $\text{range}(\mathbf{P}_{t_{\text{train}}})$ obtained from training data will actually be exact, i.e. we will have $\text{dif}(\hat{\mathbf{P}}_{t_{\text{train}}}, \mathbf{P}_{t_{\text{train}}}) = 0$. The theorem assumption that $\text{dif}(\hat{\mathbf{P}}_{t_{\text{train}}}, \mathbf{P}_{t_{\text{train}}}) \leq r_0 \zeta$ allows for the initial training data to be noisy or not exactly satisfying the model. If the training data is noisy, we need to know r_0 (in practice this is computed by thresholding to retain a certain percentage of largest eigenvalues). In this case we can let $\hat{\lambda}_{\text{train}}^-$ be the r_0 -th eigenvalue of $\frac{1}{\alpha} \sum_{t=1}^{t_{\text{train}}} \mathbf{m}_t \mathbf{m}_t'$ and $\hat{\mathbf{P}}_{t_{\text{train}}}$ be the r_0 top eigenvectors.

The following corollary is also proved when we prove the above result.

Corollary 2.2.6. *The following conclusions also hold under the assumptions of Theorem 2.2.5 or 2.2.7 with probability at least $1 - n^{-10}$*

1. The estimates of the subspace change times given by Algorithm 1 satisfy $t_j \leq \hat{t}_j \leq t_j + 2\alpha$, for $j = 1, \dots, J$;
2. The estimates of the number of new directions are correct, i.e. $\hat{r}_{j,\text{new},k} = r_{j,\text{new}}$ for $j = 1, \dots, J$ and $k = 1, \dots, K$;
3. The recovery error satisfies:

$$\|\hat{\ell}_t - \ell_t\|_2 \leq \begin{cases} 1.2(\sqrt{\zeta} + \sqrt{r_{\text{new}}}\gamma_{\text{new}}) & t \in [t_j, \hat{t}_j] \\ 1.2(1.84\sqrt{\zeta} + (0.83)^{k-1}\sqrt{r_{\text{new}}}\gamma_{\text{new}}) & t \in [\hat{t}_j + (k-1)\alpha, \hat{t}_j + k\alpha - 1], \\ & k = 1, 2, \dots, K \\ 2.4\sqrt{\zeta} & t \in [\hat{t}_j + K\alpha, t_{j+1} - 1]; \end{cases}$$

4. The subspace error satisfies,

$$\text{SE}_t \leq \begin{cases} 1 & t \in [t_j, \hat{t}_j] \\ 10^{-2}\sqrt{\zeta} + 0.83^{k-1} & t \in [\hat{t}_j + (k-1)\alpha, \hat{t}_j + k\alpha - 1], \\ & k = 1, 2, \dots, K \\ 10^{-2}\sqrt{\zeta} & t \in [\hat{t}_j + K\alpha, t_{j+1} - 1]. \end{cases}$$

2.2.5 Main Result for Online Robust PCA

Recall that in this case we assume that the observations \mathbf{m}_t satisfy $\mathbf{m}_t = \ell_t + \mathbf{x}_t$ with the support of \mathbf{x}_t , denoted \mathcal{T}_t , not known. We have the following result for Algorithm 2 given and explained in Section 2.4. This requires two extra assumptions beyond what the previous result needed. For the matrix completion problem, the set of missing entries is known, while in the robust PCA setting, the support set, \mathcal{T}_t , of the sparse outliers, \mathbf{x}_t , must be determined. We recover this using an ℓ_1 minimization step followed by thresholding. To do this correctly, we need a lower bound on the absolute values of the nonzero entries of \mathbf{x}_t . Moreover, Algorithm 2 has two extra parameters - ξ , which is the bound on the two norm of the noise seen by the ℓ_1 minimization step, and ω , which is the threshold used to recover the support of \mathbf{x}_t . These need to be set appropriately.

Theorem 2.2.7. *Consider Algorithm 2. Assume that \mathbf{m}_t satisfies (2.2) and assume everything else in Theorem 2.2.5. Also assume*

1. *The two extra algorithm parameters are set as: $\xi = \sqrt{r_{\text{new}}}\gamma_{\text{new}} + (\sqrt{r} + \sqrt{r_{\text{new}}})\sqrt{\zeta}$ and $\omega = 7\xi$*
2. *We have $x_{\min} := \min_t \min_{i: (\mathbf{x}_t)_i \neq 0} |(\mathbf{x}_t)_i| > 14\xi$*

Then with probability at least $1 - n^{-10}$,

1. *all conclusions of Theorem 2.2.5 and Corollary 2.2.6 hold;*
2. *the support set \mathcal{T}_t is exactly recovered, i.e. $\hat{\mathcal{T}}_t = \mathcal{T}_t$ for all t ;*
3. *$\|\mathbf{x}_t - \hat{\mathbf{x}}_t\|_2 = \|\boldsymbol{\ell}_t - \hat{\boldsymbol{\ell}}_t\|_2$ and $\|\boldsymbol{\ell}_t - \hat{\boldsymbol{\ell}}_t\|_2$ satisfies the bounds given in Theorem 2.2.5 and Corollary 2.2.6.*

The second assumption above can be interpreted as either a lower bound on x_{\min} , or as an upper bound on $\sqrt{r_{\text{new}}}\gamma_{\text{new}}$ in terms of x_{\min} . This latter interpretation is another “slow subspace change” condition. For the \mathbf{x}_t ’s, this result shows that their support is exactly recovered w.h.p. and its nonzero entries are accurately recovered.

2.2.6 Simple Generalizations

Consider the subspace change model, Model 2.2.2. For simplicity we put a slow increase model on the eigenvalues along the new directions for the entire period $[t_j, t_{j+1})$. However, as explained below the model, the proof of our result does not actually use this slow increase model. It only uses (3b), i.e. $\hat{\lambda}_{\text{train}}^- \leq \lambda_{\text{new}}^- \leq \lambda_{\text{new}}^+ \leq 3\hat{\lambda}_{\text{train}}^-$. Recall that λ_{new}^- and λ_{new}^+ are the minimum and maximum eigenvalues along the new directions for the first d frames after a subspace change. Thus, in the interval $[t_j + d + 1, t_{j+1})$ our proof actually does not need any constraint on $\boldsymbol{\Lambda}_{t, \text{new}}$.

With a minor modification to our proof, we can prove our result with an even weaker condition. We need (3b) to hold with λ_{new}^- being the minimum of the minimum eigenvalues of any α -frame *average* covariance matrix along the new directions over the period $[t_j, t_j + d]$, i.e.

with $\lambda_{\text{new}}^- = \min_j \min_{\tau \in [t_j, t_j + d - \alpha]} \lambda_{\min}(\frac{1}{\alpha} \sum_{t=\tau}^{\tau + \alpha - 1} \mathbf{\Lambda}_{t, \text{new}})$. For video analytics, this translates to requiring that, after a subspace change, *enough (but not necessarily all)* background frames have ‘detectable’ energy along the new directions, so that the minimum eigenvalue of the average covariance is above a threshold.

Secondly, we should point out that there is a trade off between the bound on $q_i v_i^d$, and consequently on λ_{new}^+ , in Model 2.2.2 and the bound on $\rho^2 \beta$ assumed in Model 2.2.3. Allowing a larger value of $q_i v_i^d$ (faster subspace change) will require a tighter bound on $\rho^2 \beta$ which corresponds to requiring more changes to \mathcal{T}_t . We chose the bounds $q_i (v_i)^d \leq 3$ and $\rho^2 \beta \leq .01\alpha$ for simplicity of computations. There are many other pairs that would also work. The above trade-off can be seen from the proof of Lemma 2.6.14. The proof uses Model 2.5.1 of which Model 2.2.3 is a special case. For video analytics, this means that if the background subspace changes are faster, then we also need the foreground objects to be moving more so we can ‘see’ enough of the background behind them.

Thirdly, in Model 2.2.2 we let $\mathbf{P}_t \mathbf{\Lambda}_t \mathbf{P}_t'$ be an EVD of $\mathbf{\Sigma}_t$. This automatically implies that $\mathbf{\Lambda}_t$ is diagonal. But our proof only uses the fact that $\mathbf{\Lambda}_t$ is block diagonal with blocks $\mathbf{\Lambda}_{t,*}$ and $\mathbf{\Lambda}_{t, \text{new}}$. If we relax this and we let $\mathbf{P}_t \mathbf{\Lambda}_t \mathbf{P}_t'$ be a decomposition of $\mathbf{\Sigma}_t$ where $\mathbf{\Lambda}_t$ is block diagonal as above, then our model allows the variance along *any* direction from $\text{range}(\mathbf{P}_{t_{j-1}})$ to become zero for any period of time and/or become nonzero again later. Thus, in the special case of (2.6) we can actually allow $\mathbf{P}_t = [(\mathbf{P}_{t-1} \mathbf{R}_t \setminus \mathbf{P}_{t, \text{del}}) \quad \mathbf{P}_{t, \text{new}}]$, where \mathbf{R}_t is an $r_{j-1} \times r_{j-1}$ rotation matrix and $\mathbf{P}_{t, \text{del}}$ contains the columns of $\mathbf{P}_{t-1} \mathbf{R}_t$ for which the variance is zero. This will be a special case of this generalization if $[\mathbf{P}_{t_1-1} \quad \mathbf{P}_{t_1, \text{new}} \quad \mathbf{P}_{t_2, \text{new}} \quad \dots \quad \mathbf{P}_{t_J, \text{new}}]$ is a basis matrix.

Lastly, the first condition of the theorem requires that we have accurate initial subspace knowledge. As explained below the theorem, this means that we can allow noisy training data. Moreover, notice that if we let $t_1 = t_{\text{train}} + 1$, then new background directions can enter the subspace at the same time as the first foreground object. Said another way, all we need is an accurate enough estimate of all but r_{new} directions of the initial subspace, and an assumption of small eigenvalues for sometime (d frames) along the directions for which we do not have an accurate enough estimate (or do not have an estimate).

Now consider the denseness assumption. Define the (un)denseness coefficient as follows.

Definition 2.2.8. For a basis matrix \mathbf{P} , define $\kappa_s(\mathbf{P}) := \max_{|\mathcal{T}| \leq s} \|\mathbf{I}_{\mathcal{T}}' \mathbf{P}\|_2$.

Notice that left hand side in (2.7) is $[\kappa_1(\mathbf{P})]^2$. Using the triangle inequality, it is easy to show that $\kappa_s(\mathbf{P}) \leq \sqrt{s} \kappa_1(\mathbf{P})$ [11]. Therefore, using the fact that for a basis matrix $[\mathbf{P}_1 \ \mathbf{P}_2]$, $(\kappa_s([\mathbf{P}_1 \ \mathbf{P}_2]))^2 \leq (\kappa_s(\mathbf{P}_1))^2 + (\kappa_s(\mathbf{P}_2))^2$ (see proof of the first statement of Lemma 2.C.2 in Appendix 2.C), the denseness assumptions of Theorem 2.2.7 imply that

$$\kappa_{s,*} := \kappa_{2s}(\mathbf{P}_{t_J}) \leq 0.3 \quad \text{and} \quad \kappa_{s,\text{new}} := \max_j \kappa_{2s}(\mathbf{P}_{t_j,\text{new}}) \leq 0.02. \quad (2.9)$$

The proof of Theorem 2.2.7 only uses (2.9) for the denseness assumption.

The reason for defining κ_s as above is the following lemma from [11].

Lemma 2.2.9 ([11]). For a basis matrix \mathbf{P} , $\delta_s(\mathbf{I} - \mathbf{P}\mathbf{P}') = (\kappa_s(\mathbf{P}))^2$.

2.3 Discussion

2.3.1 Discussion of the assumptions used

In the previous section, we provide two related results, one for online matrix completion (MC) and the second for online robust PCA (RPCA). The result for online RPCA can also be interpreted as a result for online sparse matrix recovery in (potentially) large but structured noise ℓ_t . Notice that our result does not require an upper bound on λ^+ (the maximum eigenvalue of $\text{Cov}(\ell_t)$ at any time) or on γ (the bound on the maximum magnitude of any entry of $\mathbf{P}_t' \ell_t$ for any time t). Both these parameters are only used to select ζ , which in turn governs the value of K and α and hence governs the required delay between subspace change times.

Our results require accurate initial subspace knowledge. As explained earlier, for video analytics, this corresponds to requiring an initial short sequence of background-only video frames whose subspace can be estimated via SVD (followed by using a singular value threshold to retain a certain number of top left singular vectors). Alternatively if an initial short sequence of the video data satisfies the assumptions required by a batch method such as PCP (for RPCA) and NNM (for MC), that can be used to estimate the low-rank part, followed by SVD to get

the column subspace. For online MC, another alternative is to use the initialization techniques of GROUSE [7] or PETRELS [8] or to use the adaptive MC idea of [10].

In Model 2.2.2, we are placing a slow increase assumption on the eigenvalues along the new directions, $\mathbf{P}_{t_j, \text{new}}$, for the interval $[t_j, t_{j+1})$. Thus after t_{j+1} , the eigenvalues along $\mathbf{P}_{t_j, \text{new}}$ can increase gradually or suddenly to any large value up to λ^+ . In fact as explained above, our proof needs the slow increase to hold only for the first d time instants after t_j , so, in fact, at any time after $t_j + d$, the eigenvalues along $\mathbf{P}_{t_j, \text{new}}$ could increase to a large value.

Model 2.2.3 on \mathcal{T}_t is a practical model for moving foreground objects in video. We should point out that this model is one special case of the general set of conditions we need (Model 2.5.1). Some other special cases of it are discussed in Section 2.9.

The model on \mathcal{T}_t (Model 2.2.3) and the denseness condition of the theorem constrain s and $s, r_0, r_{\text{new}}, J$ respectively. Model 2.2.3 requires $s \leq \rho_2 n / \alpha$ for a constant ρ_2 . Using the expression for α , it is easy to see that as long as $J \in \mathcal{O}(n)$, we have $\alpha \in \mathcal{O}(\log n)$ and so Model 2.2.3 needs $s \in \mathcal{O}(\frac{n}{\log n})$. With $s \in \mathcal{O}(\frac{n}{\log n})$, the denseness condition will hold if $r_0 \in \mathcal{O}(\log n)$, $J \in \mathcal{O}(\log n)$ and r_{new} is a constant. This is one set of sufficient conditions that we allow on the rank-sparsity product.

2.3.2 Comparison with the results for PCP and NNM

Let $\mathbf{L} := [\ell_1, \ell_2, \dots, \ell_{t_{\max}}]$ and $\mathbf{S} := [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{t_{\max}}]$. Let $r_{\text{mat}} := \text{rank}(\mathbf{L})$. Clearly $r_{\text{mat}} \leq r_0 + Jr_{\text{new}}$ and the bound is tight. Let $s_{\text{mat}} := t_{\max}s$ be a bound on the total number of missing entries of \mathbf{L} or on the support size of the outliers' matrix \mathbf{S} . In terms of r_{mat} and s_{mat} , what we need is $r_{\text{mat}} \in \mathcal{O}(\log n)$ and $s_{\text{mat}} \in \mathcal{O}(\frac{nt_{\max}}{\log n})$. This is stronger than what the PCP result from [1] or the result for NNM from [5] need (e.g., the PCP result from [1] allows $r_{\text{mat}} \in \mathcal{O}(\frac{n}{(\log n)^2})$ while allowing $s_{\text{mat}} \in \mathcal{O}(nt_{\max})$), but is similar to what the PCP results from [2, 3] need.

Other disadvantages of our result are as follows. (1) Our result needs accurate initial subspace knowledge and slow subspace change of ℓ_t . As explained earlier and in [11, Fig. 6], both of these are often practically valid for video analytics applications. Moreover, we also need the ℓ_t 's to be zero mean and mutually independent over time. Zero mean is achieved

by letting ℓ_t be the background image at time t with an empirical ‘mean background image’, computed using the training data, subtracted out. The independence assumption then models independent background variations around a common mean. As we explain in Section 2.9, this can be easily relaxed and we can get a result very similar to the current one under a first order autoregressive model on the ℓ_t ’s. (2) Moreover, Algorithms 1 and 2 need multiple algorithm parameters to be appropriately set. The PCP or NNM results need this for none [1, 5] or at most one [2, 3] algorithm parameter. (3) Thirdly, our result for online RPCA also needs a lower bound on x_{\min} while the PCP results do not need this. (4) Moreover, even with this, we can only guarantee accurate recovery of ℓ_t , while PCP or NNM guarantee exact recovery.

The advantages of our work are (1) that we analyze an online algorithm (ReProCS) that is faster and needs less storage compared with PCP or NNM. It needs to store only a few $n \times \alpha$ or $n \times r_{\text{mat}}$ matrices, thus the storage complexity is $\mathcal{O}(n \log n)$ while that for PCP or NNM is $\mathcal{O}(nt_{\max})$. In general t_{\max} can be much larger than $\log n$. (2) Moreover, we do not need any assumption on the right singular vectors of \mathbf{L} while all results for PCP or NNM do. (3) Most importantly, our results allow highly correlated changes of the set of missing entries (or outliers). From the assumption on \mathcal{T}_t , it is easy to see that we allow the number of missing entries (or outliers) per row of \mathbf{L} to be $\mathcal{O}(t_{\max})$ as long as the sets follow Model 2.2.3⁴. The PCP results from [2, 3] need this number to be $\mathcal{O}(\frac{t_{\max}}{r_{\text{mat}}})$ which is stronger. The PCP result from [1] or the NNM result [5] need an even stronger condition - they need the set $(\cup_{t=1}^{t_{\max}} \mathcal{T}_t)$ to be generated uniformly at random.

2.3.3 Other results for online RPCA and online MC

Our online RPCA result improves upon the online RPCA results from our earlier work [11] for two reasons. First, the result of [11] was a *partial result* because it required a denseness assumption on $(\mathbf{I} - \mathbf{P}_{t_j, \text{new}} \mathbf{P}_{t_j, \text{new}}') \hat{\mathbf{P}}_t$ and $(\mathbf{I} - \hat{\mathbf{P}}_{t,*} \hat{\mathbf{P}}_{t,*}' - \hat{\mathbf{P}}_{t, \text{new}} \hat{\mathbf{P}}_{t, \text{new}}') \mathbf{P}_{t_j, \text{new}}$. Here $\hat{\mathbf{P}}_{t,*}$ and $\hat{\mathbf{P}}_{t, \text{new}}$ are estimates computed by Algorithm 2. Thus, the result depended on intermediate algorithm estimates satisfying certain properties. In this work, we remove this requirement

⁴In a period of length α , the set \mathcal{T}_t can occupy index i for at most $\rho\beta$ time instants, and this pattern is allowed to repeat every α time instants. So an index can be in the support for a total of $\rho\beta \frac{t_{\max}}{\alpha}$ time instants and the model assumes $\rho\beta \leq \frac{0.01\alpha}{\rho}$ for a constant ρ .

and instead provide a *complete correctness result*. The extra assumption that we need is Model 2.2.3 on \mathcal{T}_t (or its generalization given in Model 2.5.1 later). Secondly, we provide a correctness result for a ReProCS-based algorithm that detects subspace change automatically and also estimates the rank of the new subspace automatically. The algorithm studied in [11] required knowing t_j and $r_{j,\text{new}}$ exactly for each j . Algorithms 1 and 2 in this work only require upper bounds on r_{new} , γ_{new} and J (these are needed to set the algorithm parameters - α and K for Algorithm 1, and also ξ and ω for Algorithm 2) and a small enough ζ (need bounds on r , λ^+ and γ to set this). A third minor advantage is that we also provide an algorithm and a result for online MC.

The proof of our results adapts the overall framework developed in [11]. The two important additions are: (a) Model 2.5.1 and Lemma 2.5.3 for it, and the way it is used in the proof of Lemma 2.6.23; and (b) the detection lemma (Lemma 2.6.17), the no false detection lemma (Lemma 2.6.16) and the p-PCA lemma (Lemma 2.6.18) and the lemmas used to prove these. (a) allows us to get a complete correctness result; (b) allows us to analyze an algorithm that does not use knowledge of t_j or $r_{j,\text{new}}$.

In [20], Feng et. al. propose a method for online RPCA and prove a partial result for their algorithm. The approach is to reformulate the PCP program and use this reformulation to develop a recursive algorithm that converges asymptotically to the solution of PCP as long as the basis estimate $\hat{\mathbf{P}}_t$ is full rank at each time t . Since this result assumes something about the algorithm estimates, it is also only a *partial* result.

Another recent work that uses knowledge of the initial subspace estimate and performs recovery in a piecewise batch fashion is modified-PCP [21]. However, like PCP, the result for modified PCP also needs uniformly randomly generated support sets. Its advantage is that its assumption on the rank-sparsity product is weaker than that of PCP, and hence weaker than that needed by this work. A detailed simulation comparison between modified-PCP, ReProCS and PCP demonstrating both these things is available in [21, Fig. 6].

Some other recent works that also study the online MC problem (defined differently from how we define it) include [6], Grassmanian Rank-One Update Subspace Estimation (GROUSE) [7] and Parallel Subspace Estimation and Tracking by Recursive Least Squares From Partial

Observations (PETRELS) [8]. We discuss the connection with [6] in Section 2.4. GROUSE is a first order stochastic gradient method. It uses rank-one updates to track the underlying subspace on the Grassmannian manifold. A result for its convergence to the local minimum of the cost function it optimizes is obtained in [9]. PETRELS is a second order stochastic gradient method. As explained in [8], in PETRELS, the low-dimensional subspace is tracked by minimizing a geometrically discounted sum of projection residuals on the observed entries at each time index. If missing entries are required then they can be reconstructed via least squares estimation. The subspace is updated recursively so that it is not necessary to retain historical data indefinitely. If the underlying subspace is fixed and the data stream is fully observed, then it is shown that the PETRELS estimate converges to the true subspace. In general, it always converges to the stationary point of the cost function it optimizes [8]. The advantage of PETRELS and GROUSE is that they do not need initial subspace knowledge. For our algorithms, when the initial subspace knowledge is not available or initial complete and outlier-free data is not available, we can also use the PETRELS or GROUSE ideas for initialization.

2.4 Automatic ReProCS Algorithms for Online MC and Online RPCA and Why They Work

In this section, we first introduce the automatic ReProCS based algorithm for online MC and explain why it works (this also provides the key idea why the proof of our main result would go through). Next, we do the same thing for the online RPCA algorithm. In the last two subsections (Sections 2.4.3 and 2.4.4), we explain the key insight used by our proof and give the proof outline.

2.4.1 Automatic ReProCS for Online MC (Algorithm 1)

The model on \mathbf{m}_t from (2.1) is a special case of that from (2.2) with $\mathbf{x}_t = -\mathbf{I}_{\mathcal{T}_t} \mathbf{I}_{\mathcal{T}_t}' \ell_t$ and with the support of \mathbf{x}_t , \mathcal{T}_t known [1]. Thus, we can use a simplification of the ReProCS idea for online RPCA [11] to also solve the online MC problem.

Algorithm 1 proceeds as follows. Let $\hat{\mathbf{P}}_{t-1}$ denote the basis matrix for the estimate of the subspace where ℓ_{t-1} lies. If it is an accurate estimate, because of “slow subspace change”, projecting the measurement $\mathbf{m}_t = \mathbf{x}_t + \ell_t$ onto its orthogonal complement will nullify most of ℓ_t . Specifically, we compute $\mathbf{y}_t := \Phi_t \mathbf{m}_t$ where $\Phi_t := \mathbf{I} - \hat{\mathbf{P}}_{t-1} \hat{\mathbf{P}}_{t-1}'$. Thus, \mathbf{y}_t can be rewritten as

$$\mathbf{y}_t = \Phi_t \mathbf{x}_t + \mathbf{b}_t \text{ where } \mathbf{b}_t := \Phi_t \ell_t$$

and it can be argued that $\|\mathbf{b}_t\|_2$ is small. Since the support of \mathbf{x}_t , \mathcal{T}_t , is known, we can simply recover its nonzero entries by least squares (LS) estimation, i.e. we get $\hat{\mathbf{x}}_t = \mathbf{I}_{\mathcal{T}_t}(\Phi_t)_{\mathcal{T}_t}^\dagger \mathbf{y}_t$ and then get an estimate of ℓ_t as $\hat{\ell}_t = \mathbf{m}_t - \hat{\mathbf{x}}_t$. The above approach of recovering ℓ_t is equivalent to that used by Brand in [6]; there they recover ℓ_t as an LS estimate of $\hat{\mathbf{P}}\hat{\mathbf{P}}'\ell_t \approx \ell_t$.

Let $\mathbf{e}_t := \ell_t - \hat{\ell}_t$. With the above, it is easy to see that

$$\mathbf{e}_t = \mathbf{I}_{\mathcal{T}_t}(\Phi_t)_{\mathcal{T}_t}^\dagger \mathbf{b}_t = \mathbf{I}_{\mathcal{T}_t}[(\Phi_t)_{\mathcal{T}_t}'(\Phi_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \Phi_t \ell_t.$$

Using the denseness assumption, it can be argued that the RIC of Φ_t will be small (see Lemma 2.2.9). Under the theorem’s assumptions, and conditioned on accurate recovery so far, we can bound it by 0.14. Thus, $\|(\Phi_t)_{\mathcal{T}_t}'(\Phi_t)_{\mathcal{T}_t}^{-1}\|_2 \leq 1/(1 - 0.14) < 1.2$ and so $\|\mathbf{e}_t\|_2 \leq 1.2\|\mathbf{b}_t\|_2$, i.e. it is small too (see Lemma 2.6.15).

Projection-PCA (p-PCA). The next step is to use a modification of standard PCA called projection-PCA (p-PCA), to update the subspace estimate. The reason we need p-PCA is this. Let \sum_t denote a sum over an α length time interval. In our problem, the error, \mathbf{e}_t , in the observation/estimate of ℓ_t , $\hat{\ell}_t$, is correlated with ℓ_t . Because of this, the dominant terms in the perturbation seen by standard PCA, $\frac{1}{\alpha} \sum_t \hat{\ell}_t \hat{\ell}_t' - \frac{1}{\alpha} \sum_t \ell_t \ell_t'$, are $\frac{1}{\alpha} \sum_t \ell_t \mathbf{e}_t'$ and its transpose⁵. Thus, when the condition number of $\text{Cov}(\ell_t)$ is large, it becomes difficult to argue that the perturbation will be small compared to the smallest eigenvalue of $\text{Cov}(\ell_t)$. With a large perturbation, either the $\sin \theta$ theorem [22] (that bounds the subspace error between the eigenvectors of the true and estimated sample covariance matrices) cannot be applied or it gives a useless bound.

⁵When ℓ_t and \mathbf{e}_t are uncorrelated and one of them is zero mean, it can be argued by law of large numbers that, whp, these two terms will be close to zero and $\frac{1}{\alpha} \sum_t \ell_t \ell_t'$ will be the dominant term.

Our proposed approach, projection-PCA (p-PCA) addresses the above issue as follows. At $t = t_j$, let $\mathbf{P}_* := \mathbf{P}_{t_{j-1}}$, $\mathbf{P}_{\text{new}} := \mathbf{P}_{t_j, \text{new}}$, and suppose that the subspace $\text{range}(\mathbf{P}_*)$ has been accurately recovered, i.e. we have $\hat{\mathbf{P}}_*$ so that $\text{dif}(\hat{\mathbf{P}}_*, \mathbf{P}_*) \ll 1$. Then at a time at or after $t_j + \alpha$ if we project the α previous $\hat{\ell}_t$'s perpendicular to $\hat{\mathbf{P}}_*$, we will considerably reduce the perturbation seen by the PCA step. We detect subspace change by checking if the maximum singular value of the matrix formed by these projected $\hat{\ell}_t$'s is above a threshold. Denote the time at which change is detected by \hat{t}_j . After \hat{t}_j we use SVD on K different sets of α frames of the projected $\hat{\ell}_t$'s to get improved estimates of the new subspace $\text{range}(\mathbf{P}_{\text{new}})$ in each iteration. To be precise, we get the k -th estimate, $\hat{\mathbf{P}}_{\text{new},k}$, as the left singular vectors of $(\mathbf{I} - \hat{\mathbf{P}}_* \hat{\mathbf{P}}_*') [\hat{\ell}_{\hat{t}_j + (k-1)\alpha + 1}, \dots, \hat{\ell}_{\hat{t}_j + k\alpha}]$ with singular values above a threshold. After each p-PCA step, we update $\hat{\mathbf{P}}_t$ as $\hat{\mathbf{P}}_t = [\hat{\mathbf{P}}_* \ \hat{\mathbf{P}}_{\text{new},k}]$. Finally at time $t = \hat{t}_j + K\alpha$, we update $\hat{\mathbf{P}}_*$ as $[\hat{\mathbf{P}}_* \ \hat{\mathbf{P}}_{\text{new},K}]$.

In the subspace update step, Algorithm 1 toggles between the “detect” phase and the “ppca” phase. It starts in the “detect” phase. When a subspace change is detected, i.e. at $t = \hat{t}_j$ it enters the “ppca” phase. After K iterations of p-PCA, i.e. at $t = \hat{t}_j + K\alpha + 1$, the new subspace has been accurately estimated and this time it enters the “detect” phase again.

Why p-PCA works. The reason p-PCA works is as follows. Before the first p-PCA step, i.e. for $t \in [t_j, \hat{t}_j + \alpha)$, $\hat{\mathbf{P}}_t = \hat{\mathbf{P}}_*$ and thus the noise seen by the projected sparse recovery step, $\mathbf{b}_t = \Phi \ell_t = (\mathbf{I} - \hat{\mathbf{P}}_* \hat{\mathbf{P}}_*') \ell_t$, will be the largest. Hence the error \mathbf{e}_t will also be the largest for the $\hat{\ell}_t$'s used for the first p-PCA step. However because of the projection perpendicular to $\hat{\mathbf{P}}_*$ and slow subspace change, even this error is not too large. Because of this and because \mathbf{e}_t is sparse and supported on \mathcal{T}_t and \mathcal{T}_t follows Model 2.2.3, we can argue that $\hat{\mathbf{P}}_{\text{new},1}$ is a good estimate, i.e. $\text{dif}([\hat{\mathbf{P}}_* \ \hat{\mathbf{P}}_{\text{new},1}], \mathbf{P}_{\text{new}}) \leq 0.2 < 1$. After the first p-PCA step, $\hat{\mathbf{P}}_t = [\hat{\mathbf{P}}_* \ \hat{\mathbf{P}}_{\text{new},1}]$ and this will reduce \mathbf{b}_t and hence \mathbf{e}_t for the $\hat{\ell}_t$'s in the next α frames. This and the sparseness of \mathbf{e}_t , in turn, will mean that the perturbation seen by the second p-PCA step will be smaller and so $\hat{\mathbf{P}}_{\text{new},2}$ will be a more accurate estimate of $\text{range}(\mathbf{P}_{\text{new}})$ than $\hat{\mathbf{P}}_{\text{new},1}$. This is done K times with K chosen so that $\text{dif}([\hat{\mathbf{P}}_* \ \hat{\mathbf{P}}_{\text{new},K}], \mathbf{P}_{\text{new}}) \leq r_{\text{new}} \zeta$. By the theorem assumptions, and because we can show $t_j \leq \hat{t}_j < t_j + 2\alpha$ (we explain this below), it is clear that $t_{j+1} > \hat{t}_j + K\alpha$. Thus, the new subspace added at t_j is accurately estimated before the next change time t_{j+1} .

Why \hat{t}_j are correctly detected. As explained above, we detect subspace changes by comparing the eigenvalues of $(\mathbf{I} - \hat{\mathbf{P}}_* \hat{\mathbf{P}}_*') \frac{1}{\alpha} \sum_t \hat{\ell}_t \hat{\ell}_t' (\mathbf{I} - \hat{\mathbf{P}}_* \hat{\mathbf{P}}_*')$ to a chosen threshold at every $t = u\alpha$ for $u = 1, 2, \dots, \lfloor \frac{t_{\max}}{\alpha} \rfloor$ when the algorithm is in the “detect” phase. In order to correctly detect \hat{t}_j , the algorithm first must not falsely detect new directions when none are present and it must detect subspace change within a short delay after it has occurred. The former will occur because conditioned on accurate recovery of the current subspace, $(\mathbf{I} - \hat{\mathbf{P}}_* \hat{\mathbf{P}}_*') \frac{1}{\alpha} \sum_t \hat{\ell}_t \hat{\ell}_t' (\mathbf{I} - \hat{\mathbf{P}}_* \hat{\mathbf{P}}_*')$ will have very small eigenvalues when no new directions are present. If the recovery were exact and no new directions present, this matrix would be zero. In our case, the recovery is only accurate and so we show that all eigenvalues of this matrix will be below the chosen threshold (see Lemma 2.6.16). Next consider detection of the subspace change after it has occurred. When $u = u_j := \lceil \frac{t_j}{\alpha} \rceil$, i.e. when t_j is in the interval $((u-1)\alpha + 1, u\alpha]$, not all of the ℓ_t ’s in this interval will contain new directions. Thus, depending on where in the interval t_j lies, the algorithm may or may not detect the subspace change. However, in the *next* interval, $[u_j\alpha + 1, (u_j + 1)\alpha]$, all of the ℓ_t ’s will contain new directions, and we can prove that the subspace change will be detected w.h.p. (see Lemma 2.6.17). Thus, w.h.p., either $\hat{t}_j = u_j\alpha$, or $\hat{t}_j = (u_j + 1)\alpha$. Thus, we will be able to show that $t_j \leq \hat{t}_j \leq t_j + 2\alpha$ w.h.p..

A visual description of Algorithm 1 is shown in Fig. 2.6. This figure uses Definition 2.6.4.

Algorithm 1 ReProCS for Online MC

Parameters: α, K , *Inputs:* $\hat{\mathbf{P}}_{t_{\text{train}}}, \hat{\lambda}_{t_{\text{train}}}^-, \mathbf{m}_t$ for each t , *Output:* $\hat{\ell}_t, \hat{\mathbf{P}}_t, \hat{t}_j, \hat{r}_{j,\text{new},k}$

Let $\text{thresh} = \frac{\hat{\lambda}_{t_{\text{train}}}^-}{2}$ (this is the eigenvalue threshold that will be used to detect subspace change).

Set $\hat{\mathbf{P}}_{t,*} \leftarrow \hat{\mathbf{P}}_{t_{\text{train}}}, \hat{\mathbf{P}}_{t,\text{new}} \leftarrow [\cdot], \hat{j} \leftarrow 0$, phase \leftarrow detect.

For every $t > t_{\text{train}}$, do the following:

- Compute $\mathbf{y}_t \leftarrow \Phi_t \mathbf{m}_t$ where $\Phi_t \leftarrow \mathbf{I} - \hat{\mathbf{P}}_{t-1} \hat{\mathbf{P}}_{t-1}'$
- Estimate ℓ_t : $\hat{\ell}_t \leftarrow \mathbf{m}_t - \mathbf{I}_{\mathcal{T}_t}((\Phi_t)_{\mathcal{T}_t})^\dagger \mathbf{y}_t$
- If $t \bmod \alpha \neq 0$ then $\hat{\mathbf{P}}_{t,*} \leftarrow \hat{\mathbf{P}}_{t-1,*}, \hat{\mathbf{P}}_{t,\text{new}} \leftarrow \hat{\mathbf{P}}_{t-1,\text{new}}, \hat{\mathbf{P}}_t \leftarrow [\hat{\mathbf{P}}_{t,*} \ \hat{\mathbf{P}}_{t,\text{new}}]$
- If $t \bmod \alpha = 0$ then *detection or projection PCA*
 If phase = detect then
 1. Set $u = \frac{t}{\alpha}$ and compute $\mathcal{D}_u = (\mathbf{I} - \hat{\mathbf{P}}_{u\alpha-1,*} \hat{\mathbf{P}}_{u\alpha-1,*}') [\hat{\ell}_{(u-1)\alpha+1}, \dots, \hat{\ell}_{u\alpha}]$
 2. $\hat{\mathbf{P}}_{t,*} \leftarrow \hat{\mathbf{P}}_{t-1,*}, \hat{\mathbf{P}}_{t,\text{new}} \leftarrow \hat{\mathbf{P}}_{t-1,\text{new}}, \hat{\mathbf{P}}_t \leftarrow [\hat{\mathbf{P}}_{t,*} \ \hat{\mathbf{P}}_{t,\text{new}}]$
 3. If $\lambda_{\max}(\frac{1}{\alpha} \mathcal{D}_u \mathcal{D}_u') \geq \text{thresh}$ then
 phase \leftarrow ppca, $\hat{j} \leftarrow \hat{j} + 1, k \leftarrow 0, \hat{t}_j = t$

Else (phase = ppca) then

1. Set $u = \frac{t}{\alpha}$ and compute $\mathcal{D}_u = (\mathbf{I} - \hat{\mathbf{P}}_{u\alpha-1,*} \hat{\mathbf{P}}_{u\alpha-1,*}') [\hat{\ell}_{(u-1)\alpha+1}, \dots, \hat{\ell}_{u\alpha}]$
2. $\hat{\mathbf{P}}_{t,\text{new}} \leftarrow \text{eigenvectors}(\frac{1}{\alpha} \mathcal{D}_u \mathcal{D}_u', \text{thresh}), \hat{\mathbf{P}}_{t,*} \leftarrow \hat{\mathbf{P}}_{t-1,*}, \hat{\mathbf{P}}_t \leftarrow [\hat{\mathbf{P}}_{t,*} \ \hat{\mathbf{P}}_{t,\text{new}}]$
3. $k \leftarrow k + 1$, set $\hat{r}_{j,\text{new},k} = \text{rank}(\hat{\mathbf{P}}_{t,\text{new}})$
4. If $k = K$, then
 phase \leftarrow detect, $\hat{\mathbf{P}}_{t,*} \leftarrow \hat{\mathbf{P}}_t, \hat{\mathbf{P}}_{t,\text{new}} \leftarrow [\cdot]$

$\text{eigenvectors}(\mathcal{M}, \text{thresh})$ returns a basis matrix for the span of all eigenvectors whose eigenvalue is above thresh.

2.4.2 Automatic ReProCS for online RPCA (Algorithm 2)

For online RPCA the only difference is that the support for $\mathbf{x}_t, \mathcal{T}_t$, is not known. Hence we first recover \mathbf{x}_t by ℓ_1 minimization (or any other sparse recovery method) and then estimate its support by thresholding. The rest of the steps remain the same as above.

2.4.3 Key Insight for the Proof

The argument given while explaining why p-PCA works in Section 2.4 can be formalized to show that, w.h.p., a subspace change is detected only after a change has occurred and

Algorithm 2 ReProCS for Online RPCA

Parameters: α, K, ξ, ω , *Inputs:* $\hat{\mathbf{P}}_{t_{\text{train}}}, \hat{\lambda}_{t_{\text{train}}}^-, \mathbf{m}_t$ for each t , *Output:* $\hat{\ell}_t, \hat{\mathbf{P}}_t, \hat{t}_j$

Let $\text{thresh} = \frac{\hat{\lambda}_{t_{\text{train}}}^-}{2}$. Set $\hat{\mathbf{P}}_{t,*} \leftarrow \hat{\mathbf{P}}_{t_{\text{train}}}, \hat{\mathbf{P}}_{t,\text{new}} \leftarrow [\cdot], \hat{j} \leftarrow 0$, phase \leftarrow detect.

For every $t > t_{\text{train}}$, do the following:

- Estimate \mathcal{T}_t (the support of the outlier vector \mathbf{x}_t) and \mathbf{x}_t .
 1. compute $\mathbf{y}_t \leftarrow \Phi_t \mathbf{m}_t$ where $\Phi_t \leftarrow \mathbf{I} - \hat{\mathbf{P}}_{t-1} \hat{\mathbf{P}}_{t-1}'$
 2. solve $\min_{\mathbf{x}} \|\mathbf{x}\|_1$ s.t. $\|\mathbf{y}_t - \Phi_t \mathbf{x}\|_2 \leq \xi$ and let $\hat{\mathbf{x}}_{t,\text{cs}}$ denote its solution
 3. compute $\hat{\mathcal{T}}_t = \{i : |(\hat{\mathbf{x}}_{t,\text{cs}})_i| > \omega\}$
 4. LS estimate of \mathbf{x}_t : compute $\hat{\mathbf{x}}_t = \mathbf{I}_{\hat{\mathcal{T}}_t} ((\Phi_t)_{\hat{\mathcal{T}}_t})^\dagger \mathbf{y}_t$
 - Use all steps of Algorithm 1 with $\mathcal{T}_t \leftarrow \hat{\mathcal{T}}_t$.
-

within 2α frames of the change; and that the subspace recovery error, SE_t , will decay roughly exponentially with each p-PCA iteration and become small enough after K iterations. To do this we will use the $\sin \theta$ theorem [22] (Lemma 2.6.20) followed by the matrix Hoeffding inequality [23] (Lemmas 2.7.5, 2.7.6)) to get high probability bounds on each of the terms in the subspace error bound obtained by the $\sin \theta$ theorem.

While applying the matrix Hoeffding inequality, we need to use the following key insight about the structure of $\mathbb{E}[\frac{1}{\alpha} \sum_t (\mathbf{I} - \hat{\mathbf{P}}_* \hat{\mathbf{P}}_*') \ell_t \mathbf{e}_t']$. This matrix is the dominant term in the perturbation seen by the k -th p-PCA step. Here $\mathbb{E}[\cdot]$ denotes expectation conditioned on accurate subspace recovery so far and \sum_t denotes the sum over $t \in [\hat{t}_j + (k-1)\alpha + 1, \hat{t}_j + k\alpha]$. The model on \mathcal{T}_t and the fact that \mathbf{e}_t is supported on \mathcal{T}_t can be used to show that this matrix can be written as the product of a full matrix and a block-banded matrix: for example when $\rho = 1$, the block-banded matrix will be block-diagonal, when $\rho = 2$, it will be block-tridiagonal, and so on. Also, $\mathbb{E}[\frac{1}{\alpha} \sum_t \mathbf{e}_t \mathbf{e}_t']$ will be a block banded matrix. The 2-norm of a block banded matrix is bounded by the maximum norm of any block times the number of bands in it and hence is much smaller than that of a general full matrix.

The lemma that exploits the structure of a block-banded matrix generated due to the model on \mathcal{T}_t is Lemma 2.5.3 given in Sec 2.5. This lemma is used to bound $\mathbb{E}[\frac{1}{\alpha} \sum_t (\mathbf{I} - \hat{\mathbf{P}}_* \hat{\mathbf{P}}_*') \ell_t \mathbf{e}_t']$ and $\mathbb{E}[\frac{1}{\alpha} \sum_t \mathbf{e}_t \mathbf{e}_t']$ in the proof of Lemma 2.6.23 in Section 2.7.

2.4.4 Proof Outline

We will only prove Theorem 2.2.7. Theorem 2.2.5 follows as a corollary of Theorem 2.2.7 because of the following reasons. (1) Algorithm 1 does not compute $\hat{\mathbf{x}}_t$ or its support $\hat{\mathcal{T}}_t$. For the matrix completion problem, \mathcal{T}_t is given. Thus it does not use the parameters ξ (which is the noise bound in the ℓ_1 minimization step) and ω (which is the support estimation threshold). (2) The bound on x_{\min} and the values of the parameters ξ and ω are only used in the proof of Lemma 2.6.15 to show exact support recovery, i.e. $\hat{\mathcal{T}}_t = \mathcal{T}_t$. Since for matrix completion \mathcal{T}_t is given, Theorem 2.2.5 does not need the lower bound on x_{\min} .

The proof of Theorem 2.2.7 is given in Sections 2.6 and 2.7. Before this, in the next section (Section 2.5) we give the most general model on changes in the missing/outlier entries' set \mathcal{T}_t , Model 2.5.1, and we show that Model 2.2.3 is a special case of this model. Next, we give a key lemma for sums of sparse matrices supported on rows and columns indexed by \mathcal{T}_t satisfying this model (Lemma 2.5.3).

Section 2.6 begins with defining various quantities needed for the proof. Next, we state the main lemmas used to prove the theorem, followed by the theorem's proof. There is a main lemma associated with each of the three main tasks of the algorithm: 1) accurately recovering \mathbf{x}_t and hence ℓ_t at each time t (Lemma 2.6.15), 2) detecting (subspace change) when and only when the subspace has changed, i.e. new directions have been added to the subspace (Lemmas 2.6.17 and 2.6.16), and 3) successfully estimating the dimension of the new subspace and updating its estimate by p-PCA (Lemma 2.6.18). To maintain the flow of the argument, we defer the proofs of these lemmas to the end of the section or to the appendix.

The proofs of Lemmas 2.6.21, 2.6.22, and 2.6.23 that are used together to prove Lemmas 2.6.17, 2.6.16 and 2.6.18 are rather long and are given in section 2.7. The proof of Lemma 2.6.23 uses Lemma 2.5.3 from Section 2.5.

2.5 Most General Model on Changes in \mathcal{T}_t and a Key Lemma

2.5.1 Most General Model on Changes in \mathcal{T}_t

Here we give our most general model on how \mathcal{T}_t (the set of missing entries or the support set of \mathbf{x}_t) can change. What we need to prevent is \mathcal{T}_t occupying the same indices for too many time instants in a given interval. If \mathcal{T}_t does not change ‘enough’ in a time interval of length α , we will be unable to see enough entries of a given index of ℓ_t in order to be able to accurately fill in the missing ones.

The following model quantifies ‘enough’ for our purposes. The number of time instants for which an index is part of \mathcal{T}_t is determined both by how often this set changes, and by how much it moves when it changes. The latter is parameterized by ρ which controls how much the set moves when it changes. For example $\rho = 1$ would require that distinct sets be disjoint, and $\rho = 2$ would mean that at least half of the set is displaced whenever it changes. The parameter $h^+ \in (0, 1)$ represents the maximum fraction of time for which the set remains in a given area in a time interval of length α . The smaller h^+ , the more frequently the set will need to change in order to satisfy the model. Our result requires a bound on the product $\rho^2 h^+$.

Model 2.5.1. *Let ρ be a positive integer. Split $[1, t_{\max}]$ into intervals of length α . Use $\mathcal{J}_u := [(u-1)\alpha + 1, u\alpha]$ to denote the u -th interval. For a given interval, \mathcal{J}_u , let $\mathcal{T}_{(i),u}$ for $i = 1, \dots, l_u$ be mutually disjoint subsets of $\{1, \dots, n\}$ and let $\mathcal{J}_{(i),u}, i = 1, 2, \dots, l_u$ be a partition⁶ of the interval \mathcal{J}_u so that*

$$\text{for all } t \in \mathcal{J}_{(i),u}, \mathcal{T}_t \subseteq \mathcal{T}_{(i),u} \cup \mathcal{T}_{(i+1),u} \cup \dots \cup \mathcal{T}_{(i+\rho-1),u}. \quad (2.10)$$

Define

$$h_u \left(\alpha; \{\mathcal{T}_{(i),u}\}_{i=1,\dots,l_u}, \{\mathcal{J}_{(i),u}\}_{i=1,\dots,l_u} \right) := \max_{i=1,2,\dots,l_u} |\mathcal{J}_{(i),u}| \quad (2.11)$$

and define $h_u^(\alpha)$ which takes the minimum over all choices of $\mathcal{T}_{(i),u}$ and over all choices of the*

⁶i.e. the $\mathcal{J}_{(i),u}$ ’s are mutually disjoint intervals and their union equals \mathcal{J}_u

partition $\mathcal{J}_{(i),u}$.

$$h_u^*(\alpha) := \min_{\substack{\text{all choices of mutually disjoint } \mathcal{T}_{(i),u}, i=1,2,\dots,l_u \\ \text{and all choices of mutually disjoint } \mathcal{J}_{(i),u}, i=1,2,\dots,l_u \\ \text{so that } \cup_{i=1}^{l_u} \mathcal{J}_{(i),u} = \mathcal{J}_u \text{ and (2.10) holds}}} h_u \left(\alpha; \{\mathcal{T}_{(i),u}\}_{i=1,\dots,l_u}, \{\mathcal{J}_{(i),u}\}_{i=1,\dots,l_u} \right) \quad (2.12)$$

Assume that $|\mathcal{T}_t| \leq s$ and that for all $u = 1, \dots, \lceil \frac{t_{\max}}{\alpha} \rceil$,

$$h_u^*(\alpha) \leq h^+ \alpha \quad \text{with } h^+ \leq \frac{0.01}{\rho^2}.$$

In the above model, $h_u^*(\alpha)$ provides a bound on how long \mathcal{T}_t remains in a given “area”, $\mathcal{T}_{(i),u} \cup \mathcal{T}_{(i+1),u} \cup \dots \cup \mathcal{T}_{(i+\rho-1),u}$ during the interval \mathcal{J}_u , for the best allocation of \mathcal{T}_t ’s to a given “area” and the best choice of the “areas.”

Notice that (2.10) can always be trivially satisfied by choosing $l_u = 1$, $\mathcal{T}_{(1),u} = \{1, \dots, n\}$ and $\mathcal{J}_{(1),u} = \mathcal{J}_u$, but this will give $h_u(\alpha; \cdot) = \alpha$ and hence is not a good choice. This is why we take a minimum over all choices.

Lemma 2.5.2. *Model 2.2.3 is a special case of Model 2.5.1 above with $h^+ = \frac{\beta}{\alpha}$.*

The proof is in Appendix 2.A.

Some other special cases of the above model are discussed in Section 2.9.

2.5.2 A Key Lemma that uses Model 2.5.1

Lemma 2.5.3. *Let $s_t = |\mathcal{T}_t|$. Consider a sequence of $s_t \times s_t$ symmetric positive-semidefinite matrices \mathbf{A}_t such that $\|\mathbf{A}_t\|_2 \leq \sigma^+$ for all t . Assume that the \mathcal{T}_t obey Model 2.5.1. Let $\mathbf{M} = \sum_{t \in \mathcal{J}_u} \mathbf{I}_{\mathcal{T}_t} \mathbf{A}_t \mathbf{I}_{\mathcal{T}_t}'$ be an $n \times n$ matrix (\mathbf{I} is an $n \times n$ identity matrix). Then*

$$\|\mathbf{M}\|_2 \leq \rho^2 h^+ \alpha \sigma^+ \leq 0.01 \sigma^+ \alpha$$

Proof. We will first prove the lemma for the special case when $\rho = 2$. After this, we will show how to generalize the proof when $\rho > 2$. For a given u , let $\mathcal{T}_{(i),u}$, $i = 1, 2, \dots, l_u$, and correspondingly $\mathcal{J}_{(i),u}$ denote the best choices, i.e. the choices that attain the minimum values in the definition of $h_u^*(\alpha)$.

In the rest of the proof, we remove the subscript u from l_u and from $\mathcal{T}_{(i),u}$'s for ease of notation. For simplicity of notation, we will let $\mathcal{T}_{(l+1),u} = \emptyset$.

For times $t \in \mathcal{J}_{(i),u}$, define $\mathbf{A}_{t,\text{full}}$ to be \mathbf{A}_t with rows and columns of zeros appropriately inserted so that

$$\mathbf{I}_{\mathcal{T}_t} \mathbf{A}_t \mathbf{I}_{\mathcal{T}_t}' = \mathbf{I}_{\mathcal{T}_{(i)} \cup \mathcal{T}_{(i+1)}} \mathbf{A}_{t,\text{full}} \mathbf{I}_{\mathcal{T}_{(i)} \cup \mathcal{T}_{(i+1)}}'. \quad (2.13)$$

Such an $\mathbf{A}_{t,\text{full}}$ exists because $\mathcal{T}_t \subseteq \mathcal{T}_{(i)} \cup \mathcal{T}_{(i+1)}$ for any $t \in \mathcal{J}_{(i),u}$. Notice that

$$\|\mathbf{A}_{t,\text{full}}\|_2 = \|\mathbf{A}_t\|_2, \quad (2.14)$$

because $\mathbf{A}_{t,\text{full}}$ is permutation similar to

$$\begin{bmatrix} \mathbf{A}_t & 0 \\ 0 & 0 \end{bmatrix}.$$

Since $\mathcal{T}_{(i)}$ and $\mathcal{T}_{(i+1)}$ are disjoint, we can, after permutation similarity, correspondingly partition $\mathbf{A}_{t,\text{full}}$ as

$$\begin{bmatrix} \mathbf{A}_{t,\text{full}}^{(0,0)} & \mathbf{A}_{t,\text{full}}^{(0,1)} \\ \mathbf{A}_{t,\text{full}}^{(1,0)} & \mathbf{A}_{t,\text{full}}^{(1,1)} \end{bmatrix}$$

for all $t \in \mathcal{J}_{(i),u}$.

Notice that because \mathbf{A}_t is symmetric, $\mathbf{A}_{t,\text{full}}^{(1,0)} = (\mathbf{A}_{t,\text{full}}^{(0,1)})'$. Then,

$$\begin{aligned}
\mathbf{M} &= \sum_{t \in \mathcal{J}_u} \mathbf{I}_{\mathcal{T}_t} \mathbf{A}_t \mathbf{I}_{\mathcal{T}_t}' \\
&= \sum_{i=1}^l \sum_{t \in \mathcal{J}_{(i),u}} \mathbf{I}_{\mathcal{T}_{(i)} \cup \mathcal{T}_{(i+1)}} \mathbf{A}_{t,\text{full}} \mathbf{I}_{\mathcal{T}_{(i)} \cup \mathcal{T}_{(i+1)}}' \quad \text{by (2.13)} \\
&= \sum_{i=1}^l \sum_{t \in \mathcal{J}_{(i),u}} [\mathbf{I}_{\mathcal{T}_{(i)}} \mathbf{I}_{\mathcal{T}_{(i+1)}}] \begin{bmatrix} \mathbf{A}_{t,\text{full}}^{(0,0)} & \mathbf{A}_{t,\text{full}}^{(0,1)} \\ \mathbf{A}_{t,\text{full}}^{(1,0)} & \mathbf{A}_{t,\text{full}}^{(1,1)} \end{bmatrix} \begin{bmatrix} \mathbf{I}_{\mathcal{T}_{(i)}}' \\ \mathbf{I}_{\mathcal{T}_{(i+1)}}' \end{bmatrix} \\
&= \sum_{i=1}^l \sum_{t \in \mathcal{J}_{(i),u}} \left[\mathbf{I}_{\mathcal{T}_{(i)}} \mathbf{A}_{t,\text{full}}^{(0,0)} \mathbf{I}_{\mathcal{T}_{(i)}}' + \mathbf{I}_{\mathcal{T}_{(i)}} \mathbf{A}_{t,\text{full}}^{(0,1)} \mathbf{I}_{\mathcal{T}_{(i+1)}}' + \mathbf{I}_{\mathcal{T}_{(i+1)}} \mathbf{A}_{t,\text{full}}^{(1,0)} \mathbf{I}_{\mathcal{T}_{(i)}}' + \mathbf{I}_{\mathcal{T}_{(i+1)}} \mathbf{A}_{t,\text{full}}^{(1,1)} \mathbf{I}_{\mathcal{T}_{(i+1)}}' \right] \\
&= \mathbf{I}_{\mathcal{T}_{(1)}} \left(\sum_{t \in \mathcal{J}_{(1),u}} \mathbf{A}_{t,\text{full}}^{(0,0)} \right) \mathbf{I}_{\mathcal{T}_{(1)}}' + \sum_{i=2}^l \left[\mathbf{I}_{\mathcal{T}_{(i)}} \left(\sum_{t \in \mathcal{J}_{(i-1),u}} \mathbf{A}_{t,\text{full}}^{(1,1)} + \sum_{t \in \mathcal{J}_{(i),u}} \mathbf{A}_{t,\text{full}}^{(0,0)} \right) \mathbf{I}_{\mathcal{T}_{(i)}}' \right. \\
&\quad \left. + \mathbf{I}_{\mathcal{T}_{(l)}} \left(\sum_{t \in \mathcal{J}_{(l),u}} \mathbf{A}_{t,\text{full}}^{(1,1)} \right) \mathbf{I}_{\mathcal{T}_{(l)}}' \right. \\
&\quad \left. + \sum_{i=1}^{l-1} \left[\mathbf{I}_{\mathcal{T}_{(i)}} \left(\sum_{t \in \mathcal{J}_{(i),u}} \mathbf{A}_{t,\text{full}}^{(0,1)} \right) \mathbf{I}_{\mathcal{T}_{(i+1)}}' + \mathbf{I}_{\mathcal{T}_{(i+1)}} \left(\sum_{t \in \mathcal{J}_{(i),u}} \mathbf{A}_{t,\text{full}}^{(1,0)} \right) \mathbf{I}_{\mathcal{T}_{(i)}}' \right] \right]
\end{aligned}$$

Because $\mathcal{T}_{(i)}$ and $\mathcal{T}_{(k)}$ are disjoint for $i \neq k$, \mathbf{M} has a block tridiagonal structure (by a permutation similarity if necessary):

$$\begin{bmatrix} \mathbf{B}_{(1)} & \mathbf{C}_{(1)} & 0 & 0 \\ \mathbf{C}_{(1)}' & \mathbf{B}_{(2)} & \ddots & 0 \\ 0 & \ddots & \ddots & \mathbf{C}_{(l-1)} \\ 0 & 0 & \mathbf{C}_{(l-1)}' & \mathbf{B}_{(l)} \end{bmatrix} \quad (2.15)$$

where $\mathbf{B}_{(1)} = \sum_{t \in \mathcal{J}_{(1),u}} \mathbf{A}_{t,\text{full}}^{(0,0)}$, $\mathbf{B}_{(l)} = \sum_{t \in \mathcal{J}_{(l),u}} \mathbf{A}_{t,\text{full}}^{(1,1)}$,

$$\mathbf{B}_{(i)} = \sum_{t \in \mathcal{J}_{(i-1),u}} \mathbf{A}_{t,\text{full}}^{(1,1)} + \sum_{t \in \mathcal{J}_{(i),u}} \mathbf{A}_{t,\text{full}}^{(0,0)} \quad \text{for } i = 2, 3, \dots, l \quad (2.16)$$

and

$$\mathbf{C}_{(i)} = \sum_{t \in \mathcal{J}_{(i),u}} \mathbf{A}_{t,\text{full}}^{(0,1)} \quad \text{for } i = 1, 2, \dots, l-1. \quad (2.17)$$

Now we proceed to bound $\|\mathbf{M}\|_2$.

$$\begin{aligned} \|\mathbf{M}\|_2 &= \left\| \begin{pmatrix} \mathbf{B}_{(1)} & \mathbf{C}_{(1)} & 0 & 0 \\ \mathbf{C}_{(1)}' & \ddots & \ddots & 0 \\ 0 & \ddots & \ddots & \mathbf{C}_{(l-1)} \\ 0 & 0 & \mathbf{C}_{(l-1)}' & \mathbf{B}_{(l)} \end{pmatrix} \right\|_2 \\ &\leq \left\| \begin{pmatrix} \mathbf{B}_{(1)} & 0 & 0 & 0 \\ 0 & \ddots & 0 & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & \mathbf{B}_{(l)} \end{pmatrix} \right\|_2 + \left\| \begin{pmatrix} 0 & \mathbf{C}_{(1)} & 0 & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & \mathbf{C}_{(l-1)} \\ 0 & 0 & 0 & 0 \end{pmatrix} \right\|_2 + \left\| \begin{pmatrix} 0 & 0 & 0 & 0 \\ \mathbf{C}_{(1)}' & 0 & 0 & 0 \\ 0 & \ddots & 0 & 0 \\ 0 & 0 & \mathbf{C}_{(l-1)}' & 0 \end{pmatrix} \right\|_2. \end{aligned}$$

Call the middle matrix \mathbf{C} , and observe that $\mathbf{C}\mathbf{C}'$ is block diagonal with blocks $\mathbf{C}_{(i)}\mathbf{C}_{(i)}'$.

So $\|\mathbf{C}\|_2 = \max_i \|\mathbf{C}_{(i)}\|_2$. Therefore,

$$\begin{aligned} \|\mathbf{M}\|_2 &\leq \max_i \|\mathbf{B}_{(i)}\|_2 + 2 \max_i \|\mathbf{C}_{(i)}\|_2 \\ &\leq \max_i \left\| \sum_{t \in \mathcal{J}_{(i-1),u}} \mathbf{A}_{t,\text{full}}^{(1,1)} + \sum_{t \in \mathcal{J}_{(i),u}} \mathbf{A}_{t,\text{full}}^{(0,0)} \right\|_2 + 2 \max_i \left\| \sum_{t \in \mathcal{J}_{(i),u}} \mathbf{A}_{t,\text{full}}^{(0,1)} \right\|_2 \quad \text{by (2.16) and (2.17)} \\ &\leq \max_i \left(\sum_{t \in \mathcal{J}_{(i-1),u}} \|\mathbf{A}_t\|_2 + \sum_{t \in \mathcal{J}_{(i),u}} \|\mathbf{A}_t\|_2 \right) + 2 \max_i \sum_{t \in \mathcal{J}_{(i),u}} \|\mathbf{A}_t\|_2 \quad \text{by (2.14)} \\ &\leq (\sigma^+ h_u^*(\alpha) + \sigma^+ h_u^*(\alpha)) + 2\sigma^+ h_u^*(\alpha) \leq 4\sigma^+ h^+ \alpha \end{aligned}$$

The third row used the fact that $\|\mathbf{A}_{t,\text{full}}^{(\cdot,\cdot)}\|_2 \leq \|\mathbf{A}_{t,\text{full}}\|_2 = \|\mathbf{A}_t\|_2$ for any sub-matrix of $\mathbf{A}_{t,\text{full}}$.

This finishes the proof for the $\rho = 2$ case. For this case, notice that there are 3 bands in (2.15) - the diagonal band and one band on each side of the diagonal one. When $\rho = 3$, everything will follow analogously to the above; instead of 3 bands, there will be 5 bands in the definition of \mathbf{M} and we will be able to bound its norm by

$$\begin{aligned}
\|\mathbf{M}\|_2 &\leq \max_i \left(\sum_{t \in \mathcal{J}_{(i-2),u}} \|\mathbf{A}_t\|_2 + \sum_{t \in \mathcal{J}_{(i-1),u}} \|\mathbf{A}_t\|_2 + \sum_{t \in \mathcal{J}_{(i),u}} \|\mathbf{A}_t\|_2 \right) \\
&\quad + 2 \max_i \left(\sum_{t \in \mathcal{J}_{(i-1),u}} \|\mathbf{A}_t\|_2 + \sum_{t \in \mathcal{J}_{(i),u}} \|\mathbf{A}_t\|_2 \right) \\
&\quad + 2 \max_i \sum_{t \in \mathcal{J}_{(i),u}} \|\mathbf{A}_t\|_2 \\
&\leq 3\sigma^+ h_u^*(\alpha) + 2(2\sigma^+ h_u^*(\alpha) + \sigma^+ h_u^*(\alpha)) \leq 9\sigma^+ h^+ \alpha
\end{aligned}$$

Proceeding this way, for a general ρ , there will be $1 + 2(\rho - 1) = 2\rho - 1$ bands. Any term in the central band will contain a summation of $\|\mathbf{A}_t\|_2$ over ρ sub-intervals $\mathcal{J}_{(i),u}$; any term in the first band away from the diagonal will contain this summation over $(\rho - 1)$ sub-intervals; any term in the second band away from the diagonal will contain this summation over $(\rho - 2)$ sub-intervals; and so on. Thus, we will be summing the quantity $\sigma^+ h^+ \alpha$ a total of $(\rho + 2 \sum_{i=1}^{\rho-1} i) = \rho^2$ times and so we will get $\|\mathbf{M}\|_2 \leq \rho^2 \sigma^+ h^+ \alpha$. \square

2.6 Proof of Theorem 2.2.7 and Theorem 2.2.5

As explained in Section 2.4.4, we will only prove Theorem 2.2.7. Theorem 2.2.5 follows as an easy corollary.

2.6.1 Definitions

Definition 2.6.1. Define \mathbf{e}_t to be the error made in estimating \mathbf{x}_t and ℓ_t . That is

$$\mathbf{e}_t := \hat{\mathbf{x}}_t - \mathbf{x}_t = \ell_t - \hat{\ell}_t$$

Definition 2.6.2. Define the interval

$$\mathcal{J}_u := [(u - 1)\alpha + 1, u\alpha].$$

Also define u_j to be the u such that $t_j \in \mathcal{J}_u$. That is

$$u_j := \left\lceil \frac{t_j}{\alpha} \right\rceil.$$

For the purposes of describing events before the first subspace change, let $u_0 := 0$. Also define

$$\hat{u}_j := \frac{\hat{t}_j}{\alpha}.$$

Notice from the algorithm that this will be an integer, because detection only occurs when $t \bmod \alpha = 0$.

We will show that, under appropriate conditioning, w.h.p., $\hat{u}_j = u_j$ or $\hat{u}_j = u_j + 1$.

Definition 2.6.3. Define

$$\begin{aligned} \mathbf{P}_{(j)} &:= \mathbf{P}_{t_j} \text{ for } j = 0, 1, \dots, J \\ \mathbf{P}_{(j),*} &:= \mathbf{P}_{(j-1)} = \mathbf{P}_{t_{j-1}} \text{ and } \mathbf{P}_{(j),\text{new}} := \mathbf{P}_{t_j,\text{new}} \text{ for } j = 1, \dots, J \\ \mathbf{a}_{t,*} &:= \mathbf{P}_{(j),*}' \boldsymbol{\ell}_t \text{ and } \mathbf{a}_{t,\text{new}} := \mathbf{P}_{(j),\text{new}}' \boldsymbol{\ell}_t \text{ for } t \in [t_j, t_{j+1}) \end{aligned}$$

Thus, for $t \in [t_j, t_{j+1})$, $\boldsymbol{\ell}_t$ can be written as

$$\boldsymbol{\ell}_t = [\mathbf{P}_{(j),*} \ \mathbf{P}_{(j),\text{new}}] \begin{bmatrix} \mathbf{a}_{t,*} \\ \mathbf{a}_{t,\text{new}} \end{bmatrix} = \mathbf{P}_{(j),*} \mathbf{a}_{t,*} + \mathbf{P}_{(j),\text{new}} \mathbf{a}_{t,\text{new}}$$

and $\text{Cov}(\boldsymbol{\ell}_t) = \boldsymbol{\Sigma}_t$ can be rewritten as

$$\boldsymbol{\Sigma}_t = [\mathbf{P}_{(j),*} \ \mathbf{P}_{(j),\text{new}}] \begin{bmatrix} \boldsymbol{\Lambda}_{t,*} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Lambda}_{t,\text{new}} \end{bmatrix} \begin{bmatrix} \mathbf{P}_{(j),*}' \\ \mathbf{P}_{(j),\text{new}}' \end{bmatrix}$$

Definition 2.6.4. For $j = 1, 2, \dots, J$ and $k = 1, 2, \dots, K$ define

1. $\hat{\mathbf{P}}_{(1),*} := \hat{\mathbf{P}}_{t_{\text{train}}}$ (the initial estimate) and $\hat{\mathbf{P}}_{(j),*} := \hat{\mathbf{P}}_{\hat{t}_{j-1} + K\alpha}$. If all subspace changes are correctly detected, this is the final estimate of $\mathbf{P}_{(j),*} = \mathbf{P}_{(j-1)}$.
2. $\hat{\mathbf{P}}_{(j),\text{new},0} := [\cdot]$ and $\hat{\mathbf{P}}_{(j),\text{new},k} := \hat{\mathbf{P}}_{\hat{t}_j + k\alpha, \text{new}}$. This is the k^{th} estimate of $\mathbf{P}_{(j),\text{new}}$ (again, conditioned on correct change time detection).

Notice from the algorithm that

1. $\hat{\mathbf{P}}_{t,*} = \hat{\mathbf{P}}_{(j),*}$ for all $t \in [\hat{t}_{j-1} + K\alpha, \hat{t}_j + K\alpha - 1]$
2. $\hat{\mathbf{P}}_{t,\text{new}} = \hat{\mathbf{P}}_{(j),\text{new},k}$ for all $t \in \mathcal{J}_{\hat{u}_j + (k+1)}$

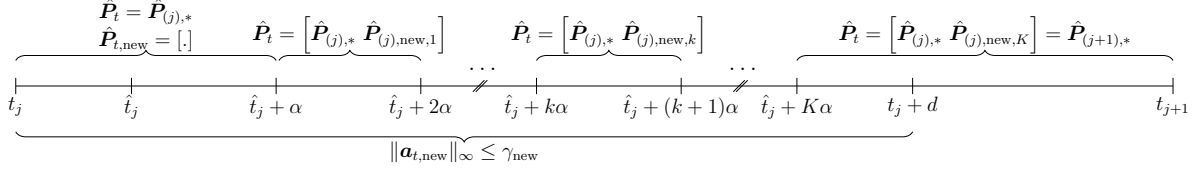


Figure 2.6: A diagram to visualize Algorithm 1 and Definition 2.6.4. The k -th projection-PCA step (at $t = \hat{t}_j + k\alpha$) computes the top left singular vectors of $(\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}') [\hat{\ell}_{\hat{t}_j + (k-1)\alpha + 1}, \hat{\ell}_{\hat{t}_j + (k-1)\alpha + 2}, \dots, \hat{\ell}_{\hat{t}_j + k\alpha}]$.

3. At all times $\hat{\mathbf{P}}_t = [\hat{\mathbf{P}}_{t,*} \ \hat{\mathbf{P}}_{t,\text{new}}]$. Thus $\hat{\mathbf{P}}_t$ and $\hat{\mathbf{P}}_{t,\text{new}}$ update at every $t = \hat{t}_j + k\alpha$, $k = 1, 2, \dots, K$, $j = 1, 2, \dots, J$ while $\hat{\mathbf{P}}_{t,*}$ updates at every $t = \hat{t}_{j-1} + K\alpha$, $j = 2, \dots, J$.
4. $\hat{\mathbf{P}}_{t-1,*} \perp \hat{\mathbf{P}}_{t,\text{new}}$ at $t = \hat{t}_j + k\alpha$ and so $\hat{\mathbf{P}}_{(j),*} \perp \hat{\mathbf{P}}_{(j),\text{new},k}$
5. $\Phi_t = (\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}' - \hat{\mathbf{P}}_{(j),\text{new},k} \hat{\mathbf{P}}_{(j),\text{new},k}')$ when $t \in \mathcal{J}_{\hat{u}_j + (k+1)}$, for $k = 1, 2, \dots, K-1$.
6. $\Phi_t = (\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}')$ when $t \in [t_j, \hat{t}_j + \alpha]$ (recall that $\hat{t}_j = \hat{u}_j \alpha$).
7. $\Phi_t = (\mathbf{I} - \hat{\mathbf{P}}_{(j+1),*} \hat{\mathbf{P}}_{(j+1),*}')$ when $t \in [\hat{t}_j + K\alpha + 1, t_{j+1} - 1]$.

Using the notation from the above definition, Figure 2.6 summarizes Algorithm 1.

Definition 2.6.5. Recall that for basis matrices \mathbf{P} and \mathbf{Q} , $\text{dif}(\mathbf{P}, \mathbf{Q}) := \|(\mathbf{I} - \mathbf{P}\mathbf{P}')\mathbf{Q}\|_2$.

Define

1. $\zeta_{j,*} := \text{dif}(\hat{\mathbf{P}}_{(j),*}, \mathbf{P}_{(j),*})$
2. $\zeta_{j,\text{new},k} := \text{dif}([\hat{\mathbf{P}}_{(j),*} \ \hat{\mathbf{P}}_{(j),\text{new},k}], \mathbf{P}_{(j),\text{new}})$

Recall $\text{SE}_t = \text{dif}(\hat{\mathbf{P}}_t, \mathbf{P}_t)$. Notice that if subspace change times are correctly detected, for $t \in \mathcal{J}_{\hat{u}_j + k}$, $\text{SE}_t \leq \zeta_{j,*} + \zeta_{j,\text{new},k-1}$ for $k = 1, 2, \dots, K$; for $t \in [t_j, \hat{t}_j + \alpha]$, $\text{SE}_t \leq 1$; and for $t \in [\hat{t}_j + K\alpha + 1, t_{j+1} - 1]$, $\text{SE}_t = \zeta_{j+1,*}$.

Definition 2.6.6. Define

1. $\zeta_{j,*}^+ := (r_0 + (j-1)r_{\text{new}})\zeta$

2. $\zeta_{j,\text{new},0}^+ := 1$, $\zeta_{j,\text{new},k}^+ := \frac{b_{\mathbf{H},k}}{b_{\mathbf{A}} - b_{\mathbf{A},\perp} - b_{\mathbf{H},k}}$ for $k = 1, 2, \dots, K$ where $b_{\mathbf{A}}$, $b_{\mathbf{A},\perp}$, and $b_{\mathbf{H},k}$ are defined in the remainder of this section. Their expressions are given by (2.21), (2.22), and (2.23).

We will show that these are high probability upper bounds on $\zeta_{j,*}$ and $\zeta_{j,\text{new},k}$ under appropriate conditioning.

As we will see later, $b_{\mathbf{A}} \approx \lambda_{\text{new}}^-$, $b_{\mathbf{A},\perp} \approx \zeta_{j,*}^{+2} \lambda^+$ and $b_{\mathbf{H},k} \approx 2\sqrt{\rho^2 h^+} \phi^+(\zeta_{j,*}^{+2} \lambda^+ + \zeta_{j,\text{new},k-1}^+ \lambda_{\text{new}}^+)$.

Here \approx means we are giving only the most dominant term for each expression. Thus,

$$\zeta_{j,\text{new},k}^+ \approx \frac{2\sqrt{\rho^2 h^+} \phi^+(\zeta_{j,\text{new},k-1}^+ \lambda_{\text{new}}^+ + \zeta_{j,*}^{+2} \lambda^+)}{\lambda_{\text{new}}^- - \zeta_{j,*}^{+2} \lambda^+ - 2\sqrt{\rho^2 h^+} \phi^+(\zeta_{j,\text{new},k-1}^+ \lambda_{\text{new}}^+ + \zeta_{j,*}^{+2} \lambda^+)}.$$

By using (2.5), the bounds on ζ from the theorem, and the bound on $\rho^2 h^+$, one can show that this decays roughly exponentially with k (see Lemma 2.6.14).

Definition 2.6.7. Define the random variable

$$X_u := \{\mathbf{a}_1, \dots, \mathbf{a}_{u\alpha}\}$$

Definition 2.6.8. Recall the definition of \mathcal{D}_u from Algorithm 1. For $j = 1, \dots, J$, $k = 1, \dots, K$, and for $a = u_j$ or $a = u_j + 1$, define the following events

- $\text{DET}_j^a := \{\hat{u}_j = a\}$
- $\text{PPCA}_{j,k}^a := \left\{ \hat{u}_j = a \text{ and } \text{rank}(\hat{\mathbf{P}}_{(j),\text{new},k}) = r_{j,\text{new}} \text{ and } \zeta_{j,\text{new},k} \leq \zeta_{j,\text{new},k}^+ \right\}$
- $\text{NODETS}_j^a := \left\{ \hat{u}_j = a \text{ and } \lambda_{\max} \left(\frac{1}{\alpha} \mathcal{D}_u \mathcal{D}_u' \right) < \text{thresh for all } u \in [\hat{u}_j + K + 1, u_{j+1} - 1] \right\}$
- $\Gamma_{0,\text{end}} := \{\zeta_{1,*} \leq r_0 \zeta\} \cap \left\{ \lambda_{\max} \left(\frac{1}{\alpha} \mathcal{D}_u \mathcal{D}_u' \right) < \text{thresh for all } u \in [1, u_1 - 1] \right\}$
- $\Gamma_{j,0}^a := \Gamma_{j-1,\text{end}} \cap \text{DET}_j^a$
- $\Gamma_{j,k}^a := \Gamma_{j,k-1}^a \cap \text{PPCA}_{j,k}^a$
- $\Gamma_{j,\text{end}} := \left(\Gamma_{j,K}^{u_j} \cap \text{NODETS}_j^{u_j} \right) \cup \left(\Gamma_{j,K}^{u_j+1} \cap \text{NODETS}_j^{u_j+1} \right)$

We misuse notation as follows. Suppose that a set Γ is a subset of all possible values that a r.v. X can take. For two r.v.s' $\{X, Y\}$, when we need to say “ $X \in \Gamma$ and Y can be anything”

we will sometimes misuse notation and just say “ $\{X, Y\} \in \Gamma$.” For example, we sometimes say $X_{u_j} \in \Gamma_{j,\text{end}}$. This means $X_{u_j-1} \in \Gamma_{j,\text{end}}$ and \mathbf{a}_t for $t \in \mathcal{J}_{u_j}$ are unconstrained.

Definition 2.6.9. Define

1. Let $\mathbf{D}_{j,\text{new}} := (\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}') \mathbf{P}_{(j),\text{new}} \stackrel{QR}{=} \mathbf{E}_{j,\text{new}} \mathbf{R}_{j,\text{new}}$ denote its reduced QR decomposition, i.e. let $\mathbf{E}_{j,\text{new}}$ be a basis matrix for $\text{range}(\mathbf{D}_{j,\text{new}})$ and let $\mathbf{R}_{j,\text{new}} = \mathbf{E}_{j,\text{new}}' \mathbf{D}_{j,\text{new}}$.
2. Let $\mathbf{E}_{j,\text{new},\perp}$ be a basis matrix for the orthogonal complement of $\text{range}(\mathbf{E}_{j,\text{new}})$. To be precise, $\mathbf{E}_{j,\text{new},\perp}$ is an $n \times (n - r_j)$ basis matrix so that $[\mathbf{E}_{j,\text{new}} \ \mathbf{E}_{j,\text{new},\perp}]$ is unitary.
3. For $u = u_j + 1$ and $u = \hat{u}_j + k$ for $k = 1, \dots, K$, define \mathbf{A}_u , $\mathbf{A}_{u,\perp}$, \mathbf{A}_u as

$$\mathbf{A}_u := \frac{1}{\alpha} \sum_{t \in \mathcal{J}_u} \mathbf{E}_{j,\text{new}}' (\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}') \ell_t \ell_t' (\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}') \mathbf{E}_{j,\text{new}}$$

$$\mathbf{A}_{u,\perp} := \frac{1}{\alpha} \sum_{t \in \mathcal{J}_u} \mathbf{E}_{j,\text{new},\perp}' (\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}') \ell_t \ell_t' (\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}') \mathbf{E}_{j,\text{new},\perp}$$

and let

$$\mathbf{A}_u := \begin{bmatrix} \mathbf{E}_{j,\text{new}} & \mathbf{E}_{j,\text{new},\perp} \end{bmatrix} \begin{bmatrix} \mathbf{A}_u & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_{u,\perp} \end{bmatrix} \begin{bmatrix} \mathbf{E}_{j,\text{new}}' \\ \mathbf{E}_{j,\text{new},\perp}' \end{bmatrix}$$

4. For $u = u_j + 1$ and $u = \hat{u}_j + k$ for $k = 1, \dots, K$, define \mathbf{M}_u and \mathbf{H}_u as

$$\mathbf{M}_u = (\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}') \left(\frac{1}{\alpha} \sum_{t \in \mathcal{J}_u} \hat{\ell}_t \hat{\ell}_t' \right) (\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}')$$

and

$$\mathbf{H}_u := \mathbf{M}_u - \mathbf{A}_u$$

Remark 2.6.10. Recall the definition of \mathcal{D}_u from Algorithm 1.

Conditioned on $\Gamma_{j-1,\text{end}}$, for $u = u_j + 1$, $\hat{\mathbf{P}}_{u\alpha-1,*} = \hat{\mathbf{P}}_{(j),*}$ (in other words all $j-1$ previous subspace changes were detected) and thus, for this value of u ,

$$\frac{1}{\alpha} \mathcal{D}_u \mathcal{D}_u' = \mathbf{M}_u.$$

In this case, \mathbf{M}_u is the matrix whose maximum eigenvalue is checked to detect subspace change.

Conditioned on $\Gamma_{j,0}^{\hat{u}_j}$, for $u = \hat{u}_j + k$, $k = 1, 2, \dots, K$, $\hat{\mathbf{P}}_{u\alpha-1,*} = \hat{\mathbf{P}}_{(j),*}$ and thus, for these values of u also,

$$\frac{1}{\alpha} \mathcal{D}_u \mathcal{D}_u' = \mathbf{M}_u.$$

In this case, \mathcal{M}_u is the matrix whose eigenvectors with eigenvalue above thresh form $\hat{\mathbf{P}}_{(j),\text{new},k}$ (see step 2 of Algorithm 1). In other words, \mathcal{M}_u has eigendecomposition

$$\mathcal{M}_u \stackrel{\text{EVD}}{=} \begin{bmatrix} \hat{\mathbf{P}}_{(j),\text{new},k} & \hat{\mathbf{P}}_{(j),\text{new},k,\perp} \end{bmatrix} \begin{bmatrix} \hat{\mathbf{\Lambda}}_u & \mathbf{0} \\ \mathbf{0} & \hat{\mathbf{\Lambda}}_{u,\perp} \end{bmatrix} \begin{bmatrix} \hat{\mathbf{P}}_{(j),\text{new},k}' \\ \hat{\mathbf{P}}_{(j),\text{new},k,\perp}' \end{bmatrix}.$$

Definition 2.6.11. Define

1. $\kappa_{s,*} := \kappa_s(\mathbf{P}_{(J)})$ and $\kappa_{s,\text{new}} := \max_j \kappa_s(\mathbf{P}_{(j),\text{new}})$.
2. $\kappa_s^+ := 0.0215$. As we will show later in Lemma 2.7.8, this upper bounds $\|\mathbf{I}_{\mathcal{T}_t}' \mathbf{D}_{j,\text{new}}\|_2$ under appropriate conditioning.
3. $\phi^+ := 1.2$. As we will show later in Lemma 2.6.15, this upper bounds $\phi_t := \|[(\Phi_t)_{\mathcal{T}_t}'(\Phi_t)_{\mathcal{T}_t}]^{-1}\|_2$ under appropriate conditioning.

Remark 2.6.12. The entire proof uses Model 2.5.1 on \mathcal{T}_t . By Lemma 2.5.2, Model 2.2.3 is a special case of it. In particular, this means that (a) Model 2.2.3 also implies $\rho^2 h^+ \leq 0.01$ and (b) Model 2.2.3 also allows us to use Lemma 2.5.3. This lemma is used in the proof of Lemma 2.6.23 in Section 2.7.

2.6.2 Five Main Lemmas for Proving Theorem 2.2.7

Fact 2.6.13. Observe that $\Gamma_{j,0}^a$ both for $a = u_j$ and $a = u_j + 1$ implies that $u_j \leq \hat{u}_j \leq u_j + 1$. Thus, in both cases, $t_j \leq \hat{t}_j \leq t_j + 2\alpha$. So with the model assumption that $d \geq (K + 2)\alpha$, we have that $\mathcal{J}_{\hat{u}_j+k} \subseteq [t_j, t_j + d]$ for $k = 1, 2, \dots, K$. This fact is needed so that we can use the “slow subspace change” inequality, (2.5), to bound the eigenvalues along the new directions, and so that we can bound $\|\mathbf{a}_{t,\text{new}}\|_\infty$ by γ_{new} .

Lemma 2.6.14. [Exponential decay of the bound on $\zeta_{j,\text{new},k}$ (similar to [11, Lemma 6.1])] Under the conditions of Theorem 2.2.7,

$$\zeta_{j,\text{new},k}^+ \leq 0.83^k + 0.84 r_{\text{new}} \zeta$$

This lemma follows by applying simple algebra on the definition and using the bounds assumed on ζ , λ_{new}^+ and $\rho^2 h^+$ in Theorem 2.2.7. A detailed proof of this lemma is given in Appendix 2.B.

Lemma 2.6.15 (Sparse Recovery Lemma (similar to [11, Lemma 6.4])). *Assume that all of the conditions of Theorem 2.2.7 hold. Recall that $\text{SE}_t = \text{dif}(\hat{\mathbf{P}}_t, \mathbf{P}_t)$.*

1. *Conditioned on $\Gamma_{j-1, \text{end}}$, for $t \in [t_j, (\hat{u}_j + 1)\alpha]$*

$$(a) \ \phi_t := \|[(\Phi_t)_{\mathcal{T}_t}'(\Phi_t)_{\mathcal{T}_t}]^{-1}\|_2 \leq \phi^+ := 1.2.$$

(b) *the support of \mathbf{x}_t is recovered exactly i.e. $\hat{\mathcal{T}}_t = \mathcal{T}_t$ and \mathbf{e}_t satisfies:*

$$\mathbf{e}_t := \hat{\mathbf{x}}_t - \mathbf{x}_t = \boldsymbol{\ell}_t - \hat{\boldsymbol{\ell}}_t = \mathbf{I}_{\mathcal{T}_t}[(\Phi_t)_{\mathcal{T}_t}'(\Phi_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \Phi_t \boldsymbol{\ell}_t. \quad (2.18)$$

(c) *Furthermore,*

$$\text{SE}_t \leq 1, \text{ and}$$

$$\|\mathbf{e}_t\|_2 \leq \phi^+ (\zeta_{j,*}^+ \sqrt{r} \gamma + \sqrt{r_{\text{new}}} \gamma_{\text{new}}) \leq 1.2 \left(\sqrt{\zeta} + \sqrt{r_{\text{new}}} \gamma_{\text{new}} \right)$$

2. *For $k = 2, 3, \dots, K$ and $\hat{u}_j = u_j$ or $\hat{u}_j = u_j + 1$, conditioned on $\Gamma_{j,k-1}^{\hat{u}_j}$, for $t \in \mathcal{J}_{\hat{u}_j+k} = [(\hat{u}_j + k - 1)\alpha + 1, (\hat{u}_j + k)\alpha]$, the first two conclusions above hold. That is, $\phi_t \leq \phi^+$ and \mathbf{e}_t satisfies (2.18). Furthermore,*

$$\text{SE}_t \leq \zeta_{j,*}^+ + \zeta_{j,\text{new},k-1}^+, \text{ and}$$

$$\|\mathbf{e}_t\|_2 \leq \phi^+ (\zeta_{j,*}^+ \sqrt{r} \gamma + \zeta_{j,\text{new},k-1}^+ \sqrt{r_{\text{new}}} \gamma_{\text{new}}) \leq 1.2 \left(1.84 \sqrt{\zeta} + (0.83)^{k-1} \sqrt{r_{\text{new}}} \gamma_{\text{new}} \right)$$

3. *For $\hat{u}_j = u_j$ or $\hat{u}_j = u_j + 1$, conditioned on $\Gamma_{j,K}^{\hat{u}_j}$, for $t \in [(\hat{u}_j + K)\alpha + 1, t_{j+1} - 1]$, the first two conclusions above hold ($\phi_t \leq \phi^+$ and \mathbf{e}_t satisfies (2.18)). Furthermore,*

$$\text{SE}_t \leq \zeta_{j+1,*}^+, \text{ and}$$

$$\|\mathbf{e}_t\|_2 \leq \phi^+ \zeta_{j+1,*}^+ \sqrt{r} \gamma \leq 1.2 \sqrt{\zeta}$$

Notice that cases 1) and 3) of the above lemma occur when the algorithm is in the detection phase, while during the intervals for case 2) the algorithm is performing projection-PCA. In case 1) new directions have been added but not estimated, so the error is larger. In case 2), the error is decaying exponentially with each estimation step. Finally, case 3) occurs after the new directions have been successfully estimated and contains the tightest error bounds.

The proof is given in Appendix 2.C.

Lemma 2.6.16 (No false detection of subspace changes).

1. The event $\Gamma_{j,K}^{\hat{u}_j}$ and so also the event $\Gamma_{j,\text{end}}$ imply that $\zeta_{j+1,*} \leq \zeta_{j+1,*}^+$.
2. $\mathbb{P}(\text{NODETS}_j^a \mid \Gamma_{j,K}^a) = 1$ for $a = u_j$ or $a = u_j + 1$.

Lemma 2.6.17 (Subspace change detected within 2α frames). For $j = 1, \dots, J$,

$$\mathbb{P}(\text{DET}^{u_j+1} \mid \Gamma_{j-1,\text{end}}, \overline{\text{DET}^{u_j}}) \geq p_{\text{det},1} := 1 - p_{\mathbf{A}} - p_{\mathbf{H}}.$$

The definitions of $p_{\mathbf{A}}$ and $p_{\mathbf{H}}$ can be found in the remainder of this section.

Lemma 2.6.18 (k -th iteration of pPCA works well).

$$\mathbb{P}(\Gamma_{j,k}^a \mid \Gamma_{j,k-1}^a) = \mathbb{P}(\text{PPCA}_{j,k}^a \mid \Gamma_{j,k-1}^a) \geq p_{\text{ppca}} := 1 - p_{\mathbf{A}} - p_{\mathbf{A},\perp} - p_{\mathbf{H}}$$

for $a = u_j$ or $a = u_j + 1$. The definitions of $p_{\mathbf{A}}$, $p_{\mathbf{A},\perp}$, and $p_{\mathbf{H}}$ can be found in the remainder of this section.

The above lemma says that, conditioned on $k - 1$ previous successful p-PCA steps and on accurate recovery of $P_{(j-1),*}$, the probability of correctly estimating $r_{j,\text{new}}$ and of a successful k^{th} projection PCA step is lower bounded by p_{ppca} . This is true whether the new directions are detected at u_j or at $u_j + 1$.

2.6.3 Proof of Theorem 2.2.7

Corollary 2.6.19. The above lemmas imply that $\mathbb{P}(\Gamma_{j,\text{end}} \mid \Gamma_{j-1,\text{end}}) \geq p_{\text{det},1} \cdot (p_{\text{ppca}})^K$.

Proof. Let

$$p_{\text{det},0} := \mathbb{P}(\text{DET}^{u_j} \mid \Gamma_{j-1,\text{end}}).$$

From the above lemmas, we get that

$$\begin{aligned}
\mathbb{P}(\Gamma_{j,\text{end}} \mid \Gamma_{j-1,\text{end}}) &= \mathbb{P}\left((\text{DET}^{u_j} \cap \text{PPCA}_{j,1}^{u_j} \cap \cdots \cap \text{PPCA}_{j,K}^{u_j} \cap \text{NODETS}_j^{u_j}) \cup \right. \\
&\quad \left. (\overline{\text{DET}^{u_j}} \cap \text{DET}^{u_j+1} \cap \text{PPCA}_{j,k}^{u_j+1} \cap \cdots \cap \text{PPCA}_{j,K}^{u_j+1} \cap \text{NODETS}_j^{u_j+1}) \mid \Gamma_{j-1,\text{end}}\right) \\
&= \mathbb{P}\left(\text{DET}^{u_j} \cap \text{PPCA}_{j,1}^{u_j} \cap \cdots \cap \text{PPCA}_{j,K}^{u_j} \mid \Gamma_{j-1,\text{end}}\right) \\
&\quad + \mathbb{P}\left(\overline{\text{DET}^{u_j}} \cap \text{DET}^{u_j+1} \cap \text{PPCA}_{j,k}^{u_j+1} \cap \cdots \cap \text{PPCA}_{j,K}^{u_j+1} \mid \Gamma_{j-1,\text{end}}\right) \\
&\geq p_{\text{det},0} \cdot (p_{\text{ppca}})^K + (1 - p_{\text{det},0}) \cdot p_{\text{det},1} \cdot (p_{\text{ppca}})^K \\
&\geq p_{\text{det},0} \cdot p_{\text{det},1} \cdot (p_{\text{ppca}})^K + (1 - p_{\text{det},0}) \cdot p_{\text{det},1} \cdot (p_{\text{ppca}})^K \\
&= p_{\text{det},1} \cdot (p_{\text{ppca}})^K.
\end{aligned}$$

□

Proof of Theorem 2.2.7. Theorem 2.2.7 follows from Corollary 2.6.19 and the assumed lower bound on α . Notice that by Lemma 2.6.14, the choice of K , and Lemma 2.6.15, the event $\Gamma_{J,\text{end}}$ will imply all conclusions of the theorem.

By the first assumption (accurate initial subspace knowledge) and the argument used to prove Lemma 2.6.16, we get that $\mathbb{P}(\Gamma_{0,\text{end}}) = 1$. By the chain rule,

$$\mathbb{P}(\Gamma_{J,\text{end}}) = \prod_{j=1}^J \mathbb{P}(\Gamma_{j,\text{end}} \mid \Gamma_{j-1,\text{end}}, \Gamma_{j-2,\text{end}}, \dots, \Gamma_{1,\text{end}}, \Gamma_{0,\text{end}}).$$

Because $\Gamma_{j-1,\text{end}} \subseteq \Gamma_{j-2,\text{end}} \subseteq \cdots \subseteq \Gamma_{1,\text{end}} \subseteq \Gamma_{0,\text{end}}$, we get

$$\begin{aligned}
\mathbb{P}(\Gamma_{J,\text{end}}) &= \prod_{j=1}^J \mathbb{P}(\Gamma_{j,\text{end}} \mid \Gamma_{j-1,\text{end}}) \\
&\geq \prod_{j=1}^J p_{\text{det},1} \cdot (p_{\text{ppca}})^K = (p_{\text{det},1} \cdot (p_{\text{ppca}})^K)^J \\
&\geq 1 - n^{-10}
\end{aligned}$$

The last line is by the lower bound on α assumed in the theorem and the fact that $p_{\text{det},1} \geq p_{\text{ppca}}$.

□

2.6.4 Key Lemmas for Proving of Lemmas 2.6.16, 2.6.17, and 2.6.18

Before proving the lemmas from the preceding subsection, we introduce several lemmas which will be used in the proofs.

The following lemma follows from the $\sin \theta$ theorem [22] and Weyl's theorem. It is taken from [11].

Lemma 2.6.20 ([11], Lemma 6.9). *At $u = \hat{u}_j + k$, if $\text{rank}(\hat{\mathbf{P}}_{(j),\text{new},k}) = r_{j,\text{new}}$, and if $\lambda_{\min}(\mathbf{A}_u) - \|\mathbf{A}_{u,\perp}\|_2 - \|\mathbf{H}_u\|_2 > 0$, then*

$$\zeta_{j,\text{new},k} \leq \frac{\|\mathbf{H}_u\|_2}{\lambda_{\min}(\mathbf{A}_u) - \|\mathbf{A}_{u,\perp}\|_2 - \|\mathbf{H}_u\|_2} \quad (2.19)$$

The next three lemmas each assert a high probability bound for one of the terms in (2.19). In the following lemmas, let

$$\epsilon = \frac{r_{\text{new}} \zeta \hat{\lambda}_{\text{train}}^-}{100}. \quad (2.20)$$

Let $p_{\mathbf{A}} := r_{\text{new}} \exp\left(\frac{-\alpha \zeta^2 (\hat{\lambda}_{\text{train}}^-)^2}{8 \cdot 100^2 \cdot \gamma_{\text{new}}^4}\right) + r_{\text{new}} \exp\left(\frac{-\alpha r_{\text{new}}^2 \zeta^2 (\hat{\lambda}_{\text{train}}^-)^2}{8 \cdot 100^2 \cdot 4^2}\right)$ and

$$b_{\mathbf{A}} := (1 - (\zeta_{j,*}^+)^2) \lambda_{\text{new}}^- - 2\epsilon. \quad (2.21)$$

Lemma 2.6.21. *For $k = 1, \dots, K$,*

$$\mathbb{P}(\lambda_{\min}(\mathbf{A}_{\hat{u}_j+k}) \geq b_{\mathbf{A}} \mid X_{\hat{u}_j+k-1}) \geq 1 - p_{\mathbf{A}}$$

for all $X_{\hat{u}_j+k-1} \in \Gamma_{j,k-1}^{\hat{u}_j}$ with $\hat{u}_j = u_j$ or $\hat{u}_j = u_j + 1$.

The same bound holds for $\lambda_{\min}(\mathbf{A}_{u_j+1})$ when we condition on $X_{u_j} \in \Gamma_{j-1,\text{end}}$.

Let $p_{\mathbf{A},\perp} := (n - r_{\text{new}}) \exp\left(\frac{-\alpha r_{\text{new}}^2 \zeta^2 (\hat{\lambda}_{\text{train}}^-)^2}{8 \cdot 100^2}\right)$

and

$$b_{\mathbf{A},\perp} := (\zeta_{j,*}^+)^2 \lambda^+ + \epsilon. \quad (2.22)$$

Lemma 2.6.22. *For $k = 1, \dots, K$,*

$$\mathbb{P}(\lambda_{\max}(\mathbf{A}_{\hat{u}_j+k,\perp}) \leq b_{\mathbf{A},\perp} \mid X_{\hat{u}_j+k-1}) \geq 1 - p_{\mathbf{A},\perp}$$

for all $X_{\hat{u}_j+k-1} \in \Gamma_{j,k-1}^{\hat{u}_j}$ with $\hat{u}_j = u_j$ or $\hat{u}_j = u_j + 1$.

The same bound holds for $\lambda_{\max}(\mathbf{A}_{u_j+1,\perp})$ when we condition on $X_{u_j} \in \Gamma_{j-1,\text{end}}$.

Let

$$\begin{aligned} p_{\mathbf{H}} := & n \exp \left(\frac{-\alpha r_{\text{new}}^2 \zeta^2 (\hat{\lambda}_{\text{train}}^-)^2}{32 \cdot 100^2 (\phi^+)^2 (\sqrt{\zeta} + \sqrt{r_{\text{new}} \gamma_{\text{new}}})^4} \right) \\ & + n \exp \left(\frac{-\alpha r_{\text{new}}^2 \zeta^2 (\hat{\lambda}_{\text{train}}^-)^2}{8 \cdot 100^2 \left(\phi^+ (\sqrt{\zeta} + \sqrt{r_{\text{new}} \gamma_{\text{new}}}) \right)^4} \right) + \\ & n \exp \left(\frac{-\alpha r_{\text{new}}^2 \zeta^2 (\hat{\lambda}_{\text{train}}^-)^2}{32 \cdot 100^2 (\zeta + \sqrt{\zeta} \sqrt{r_{\text{new}} \gamma_{\text{new}}})^2} \right). \end{aligned}$$

and

$$b_{\mathbf{H},k} := 2b_{\ell e,k} + b_{ee,k} + 2b_{\mathbf{F}} \quad (2.23)$$

where

$$\begin{aligned} b_{\ell e,k} &:= \begin{cases} \phi^+ (\sqrt{\rho^2 h^+} (\zeta_{j,*}^+)^2 \lambda^+ + \kappa_s^+ \lambda_{\text{new}}^+) + \epsilon & k = 1 \\ \left[(\zeta_{j,*}^+)^2 \lambda^+ + \zeta_{j,\text{new},k-1}^+ \lambda_{\text{new}}^+ \right] (\sqrt{\rho^2 h^+} \phi^+) + \epsilon & k \geq 2 \end{cases} \\ b_{ee,k} &:= \begin{cases} \rho^2 h^+ (\phi^+)^2 ((\zeta_{j,*}^+)^2 \lambda^+ + (\kappa_s^+)^2 \lambda_{\text{new}}^+) + \epsilon & k = 1 \\ \rho^2 h^+ (\phi^+)^2 ((\zeta_{j,*}^+)^2 (\lambda^+) + (\zeta_{j,\text{new},k-1}^+)^2 (\lambda_{\text{new}}^+)) + \epsilon & k \geq 2 \end{cases} \end{aligned}$$

and

$$b_{\mathbf{F}} := (\zeta_{j,*}^+)^2 \lambda^+ + \epsilon.$$

Lemma 2.6.23. *For $k = 1, \dots, K$,*

$$\mathbb{P} \left(\|\mathbf{H}_{\hat{u}_j+k}\|_2 \leq b_{\mathbf{H},k} \mid X_{\hat{u}_j+k-1} \right) \geq 1 - p_{\mathbf{H}} \quad (2.24)$$

for all $X_{\hat{u}_j+k-1} \in \Gamma_{j,k-1}^{\hat{u}_j}$ with $\hat{u}_j = u_j$ or $\hat{u}_j = u_j + 1$

The same bound ($k = 1$ case), i.e. $\|\mathbf{H}_{u_j+1}\|_2 \leq b_{\mathbf{H},1}$, also holds with the same probability when we condition on $X_{u_j} \in \Gamma_{j-1,\text{end}}$.

The above lemmas are proved in the next section (Section 2.7). The proofs use Fact 2.6.13.

2.6.5 Proofs of Lemmas 2.6.16, 2.6.17, and 2.6.18

Proof of Lemma 2.6.16. Recall that $\Gamma_{j,\text{end}} := \left(\Gamma_{j,K}^{u_j} \cap \text{NODETS}_j^{u_j} \right) \cup \left(\Gamma_{j,K}^{u_j+1} \cap \text{NODETS}_j^{u_j+1} \right)$.

1. By the definition of $\Gamma_{j,K}^{\hat{u}_j}$, both for $\hat{u}_j = u_j$ and $\hat{u}_j = u_j + 1$, $\zeta_{j,*} \leq \zeta_{j,*}^+ = (r_0 + (j-1)r_{\text{new}})\zeta$ and $\zeta_{j,K} \leq \zeta_{j,\text{new},K}^+$. Lemma 2.6.14 and the choice of K imply that $\zeta_{j,\text{new},K}^+ \leq r_{\text{new}}\zeta$. Thus, $\zeta_{j+1,*} \leq \zeta_{j,*} + \zeta_{j,\text{new},K} \leq \zeta_{j+1,*}^+ = (r_0 + jr_{\text{new}})\zeta$.
2. $\mathbb{P}(\text{NODETS}_j^{\hat{u}_j} \mid \Gamma_{j,K}^{\hat{u}_j}) = \mathbb{P}\left(\lambda_{\max}\left(\frac{1}{\alpha}\mathcal{D}_u\mathcal{D}_u\right) < \text{thresh for all } u \in [\hat{u}_j + K + 1, u_{j+1} - 1] \mid \Gamma_{j,K}^{\hat{u}_j}\right)$ for $\hat{u}_j = u_j$ or $\hat{u}_j = u_j + 1$.

As shown in 1), $\Gamma_{j,K}^{\hat{u}_j}$ implies that $\text{dif}(\hat{\mathbf{P}}_{(j+1),*}, \mathbf{P}_{(j+1),*}) \leq \zeta_{j+1,*}^+ = (r_0 + jr_{\text{new}})\zeta$. Recall that $\mathbf{P}_{(j+1),*} = \mathbf{P}_{(j)}$. Also, for $u \in [\hat{u}_j + K + 1, u_{j+1} - 1]$, $\hat{\mathbf{P}}_{u\alpha-1,*} = \hat{\mathbf{P}}_{(j+1),*}$. Also, for all $t \in \mathcal{J}_u$ for these u 's, $\boldsymbol{\ell}_t = \mathbf{P}_{(j)}\mathbf{a}_t = \mathbf{P}_{(j+1),*}\mathbf{a}_t$. Therefore,

$$\begin{aligned}
\lambda_{\max}\left(\frac{1}{\alpha}\mathcal{D}_u\mathcal{D}_u\right) &= \lambda_{\max}\left(\frac{1}{\alpha}\sum_{t \in \mathcal{J}_u}(\mathbf{I} - \hat{\mathbf{P}}_{u\alpha-1,*}\hat{\mathbf{P}}_{u\alpha-1,*}')\hat{\boldsymbol{\ell}}_t\hat{\boldsymbol{\ell}}_t'(\mathbf{I} - \hat{\mathbf{P}}_{u\alpha-1,*}\hat{\mathbf{P}}_{u\alpha-1,*}')\right) \\
&= \lambda_{\max}\left(\frac{1}{\alpha}\sum_{t \in \mathcal{J}_u}(\mathbf{I} - \hat{\mathbf{P}}_{(j+1),*}\hat{\mathbf{P}}_{(j+1),*}')(\mathbf{P}_{(j)}\mathbf{a}_t - \mathbf{e}_t)\right. \\
&\quad \left.(\mathbf{P}_{(j)}\mathbf{a}_t - \mathbf{e}_t)'(\mathbf{I} - \hat{\mathbf{P}}_{(j+1),*}\hat{\mathbf{P}}_{(j+1),*}')\right) \\
&\leq (\zeta_{j+1,*}^+)^2 r\gamma^2 + 2\phi^+(\zeta_{j+1,*}^+)^2 r\gamma^2 + (\phi^+)^2 (\zeta_{j+1,*}^+)^2 r\gamma^2 \\
&\leq 4(\phi^+)^2 \zeta \hat{\lambda}_{\text{train}}^- \leq \frac{\hat{\lambda}_{\text{train}}^-}{2}.
\end{aligned}$$

The bound on \mathbf{e}_t comes from Lemma 2.6.15. The penultimate inequality uses the bound $\zeta \leq \frac{\hat{\lambda}_{\text{train}}^-}{r^3\gamma^2}$ assumed in Theorem 2.2.7. \square

The next two proofs follow using the following two facts and the four lemmas from the previous subsection.

Fact 2.6.24. *For an event \mathcal{E} and random variable X , $\mathbb{P}(\mathcal{E}|X) \geq p$ for all $X \in \mathcal{C}$ implies that $\mathbb{P}(\mathcal{E}|X \in \mathcal{C}) \geq p$.*

Fact 2.6.25. *Using the bounds on ζ and on $\rho^2 h^+$ and using (3b), we get*

$$\begin{aligned}
b_{\mathbf{A}} &\geq 0.94\lambda_{\text{new}}^- \geq 0.94\hat{\lambda}_{\text{train}}^- \\
b_{\mathbf{A},\perp} &\leq 0.011\hat{\lambda}_{\text{train}}^- \\
b_{\boldsymbol{\mathcal{H}},k} &\leq 0.24\hat{\lambda}_{\text{train}}^-.
\end{aligned}$$

Thus, $b_{\mathbf{A}} - b_{\boldsymbol{\mathcal{H}},k} \geq 0.5\hat{\lambda}_{\text{train}}^- = \text{thresh}$ and $b_{\mathbf{A},\perp} + b_{\boldsymbol{\mathcal{H}},k} < 0.25\hat{\lambda}_{\text{train}}^- < \text{thresh}$.

Proof of Lemma 2.6.17. We will prove that $\mathbb{P}(\text{DET}^{u_j+1} \mid X_{u_j}) > p_{\text{det},1}$ for all $X_{u_j} \in \Gamma_{j-1,\text{end}}$. In particular, this will imply that $\mathbb{P}(\text{DET}^{u_j+1} \mid X_{u_j}) > p_{\text{det},1}$ for all $X_{u_j} \in \Gamma_{j-1,\text{end}} \cap \overline{\text{DET}^{u_j}}$ and so we can conclude that $\mathbb{P}(\text{DET}^{u_j+1} \mid \Gamma_{j-1,\text{end}}, \overline{\text{DET}^{u_j}}) > p_{\text{det},1}$.

Recall that $\mathcal{M}_u = \frac{1}{\alpha} \mathcal{D}_u \mathcal{D}_u'$, and observe that

$$\mathbb{P}(\text{DET}^{u_j+1} \mid X_{u_j}) = \mathbb{P}(\lambda_{\max}(\mathcal{M}_{u_j+1}) > \text{thresh} \mid X_{u_j})$$

By Weyl's Theorem

$$\begin{aligned} \lambda_{\max}(\mathcal{M}_{u_j+1}) &\geq \lambda_{\max}(\mathcal{A}_{u_j+1}) + \lambda_{\min}(\mathcal{H}_{u_j+1}) \\ &\geq \lambda_{\max}(\mathcal{A}_{u_j+1}) - \|\mathcal{H}_{u_j+1}\|_2 \\ &\geq \lambda_{\min}(\mathcal{A}_{u_j+1}) - \|\mathcal{H}_{u_j+1}\|_2 \end{aligned}$$

When $X_{u_j} \in \Gamma_{j-1,\text{end}}$, Lemmas 2.6.21 and 2.6.23 applied with ϵ given by (2.20) show that $\lambda_{\min}(\mathcal{A}_{u_j+1}) \geq b_{\mathcal{A}}$ and $\|\mathcal{H}_{u_j+1}\|_2 \leq b_{\mathcal{H},1}$ with probability at least $1 - p_{\mathcal{A}} - p_{\mathcal{H}} = p_{\text{det},1}$. Using Fact 2.6.25, $b_{\mathcal{A}} - b_{\mathcal{H},1} \geq \text{thresh}$ and so the lemma follows. \square

Proof of Lemma 2.6.18. To prove this Lemma we need to show two things. First, conditioned on $\Gamma_{j,k-1}^{\hat{u}_j}$, the k^{th} estimate of the number of new directions is correct. That is: $\hat{r}_{j,\text{new},k} = r_{j,\text{new}}$. Second, we must show $\zeta_{j,\text{new},k} \leq \zeta_{j,\text{new},k}^+$, again conditioned on $\Gamma_{j,k-1}^{\hat{u}_j}$.

Notice that $\hat{r}_{j,\text{new},k} = \text{rank}(\hat{\mathcal{P}}_{(j),\text{new},k})$. To show that $\text{rank}(\hat{\mathcal{P}}_{(j),\text{new},k}) = r_{j,\text{new}}$, we need to show that for $u = \hat{u}_j + k$, $k = 1, \dots, K$, $\lambda_{r_{j,\text{new}}}(\mathcal{M}_u) > \text{thresh}$ and $\lambda_{r_{j,\text{new}}+1}(\mathcal{M}_u) < \text{thresh}$. To do this we proceed similarly to above.

Observe that, $\mathcal{M}_u = \mathcal{A}_u + \mathcal{H}_u$. By Fact 2.6.25, $b_{\mathcal{A}} > b_{\mathcal{A},\perp}$. Combining this with Lemmas 2.6.21 and 2.6.22 gives, $\lambda_{\min}(\mathcal{A}_u) > \lambda_{\max}(\mathcal{A}_{u,\perp})$ with probability at least $1 - p_{\mathcal{A}} - p_{\mathcal{A},\perp}$ under the appropriate conditioning (conditioned on $\Gamma_{j,k-1}^{\hat{u}_j}$). Since \mathcal{A}_u is of size $r_{j,\text{new}} \times r_{j,\text{new}}$, this means that $\lambda_{r_{j,\text{new}}}(\mathcal{A}_u) = \lambda_{\min}(\mathcal{A}_u)$ and $\lambda_{r_{j,\text{new}}+1}(\mathcal{A}_u) = \lambda_{\max}(\mathcal{A}_{u,\perp})$. Using this and Weyl's Theorem,

$$\begin{aligned} \lambda_{r_{j,\text{new}}}(\mathcal{M}_u) &\geq \lambda_{r_{j,\text{new}}}(\mathcal{A}_u) + \lambda_{\min}(\mathcal{H}_u) \\ &\geq \lambda_{r_{j,\text{new}}}(\mathcal{A}_u) - \|\mathcal{H}_u\|_2 \\ &= \lambda_{\min}(\mathcal{A}_u) - \|\mathcal{H}_u\|_2 \end{aligned}$$

and

$$\begin{aligned}
\lambda_{r_{j,\text{new}}+1}(\mathcal{M}_u) &\leq \lambda_{r_{j,\text{new}}+1}(\mathcal{A}_u) + \lambda_{\max}(\mathcal{H}_u) \\
&\leq \lambda_{r_{j,\text{new}}+1}(\mathcal{A}_u) + \|\mathcal{H}_u\|_2 \\
&= \lambda_{\max}(\mathbf{A}_{u,\perp}) + \|\mathcal{H}_u\|_2
\end{aligned}$$

with probability at least $1 - p_{\mathbf{A}} - p_{\mathbf{A},\perp}$ under the appropriate conditioning. Using Lemmas 2.6.21, 2.6.22, and 2.6.23 applied with ϵ given by (2.20) and Fact 2.6.25, we can conclude that with probability greater than p_{ppca} , $\lambda_{r_{j,\text{new}}}(\mathcal{M}_u) > b_{\mathbf{A}} - b_{\mathcal{H},k} \geq \text{thresh}$ and $\lambda_{r_{j,\text{new}}+1}(\mathcal{M}_u) < b_{\mathbf{A},\perp} + b_{\mathcal{H},k} < \text{thresh}$. Therefore $\text{rank}(\hat{\mathbf{P}}_{(j),\text{new},k}) = r_{j,\text{new}}$ with probability greater than p_{ppca} under the appropriate conditioning.

To show that $\zeta_{j,\text{new},k} \leq \zeta_{j,\text{new},k}^+$, we also use Lemmas 2.6.21, 2.6.22, and 2.6.23 applied with ϵ given by (2.20). Using $\text{rank}(\hat{\mathbf{P}}_{(j),\text{new},k}) = r_{j,\text{new}}$ and applying Lemma 2.6.20 with these bounds gives the desired result. \square

2.7 Proofs of Lemmas 2.6.21, 2.6.22, and 2.6.23

2.7.1 Some definitions, remarks and facts

Definition 2.7.1. Define the following for $k = 0, 1, \dots, K$. Recall that $\hat{\mathbf{P}}_{(j),\text{new},0} = [\cdot]$.

1. $\mathbf{D}_{j,\text{new},k} := (\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}' - \hat{\mathbf{P}}_{(j),\text{new},k} \hat{\mathbf{P}}_{(j),\text{new},k}') \mathbf{P}_{(j),\text{new}}$. Thus $\mathbf{D}_{j,\text{new}} = \mathbf{D}_{j,\text{new},0}$.
2. $\mathbf{D}_{j,*,k} := (\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}' - \hat{\mathbf{P}}_{(j),\text{new},k} \hat{\mathbf{P}}_{(j),\text{new},k}') \mathbf{P}_{(j),*}$ and $\mathbf{D}_{j,*} := \mathbf{D}_{j,*,0}$.
3. Recall that $\zeta_{j,\text{new},0} = \|\mathbf{D}_{j,\text{new}}\|_2$, $\zeta_{j,\text{new},k} = \|\mathbf{D}_{j,\text{new},k}\|_2$, $\zeta_{j,*} = \|\mathbf{D}_{j,*}\|_2$. Also, clearly, $\|\mathbf{D}_{j,*,k}\|_2 \leq \|\mathbf{D}_{j,*}\|_2 \leq \zeta_{j,*}$.

Definition 2.7.2. For ease of notation, define

$$\tilde{\ell}_t := (\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}') \ell_t$$

Remark 2.7.3. In the rest of this section, for ease of notation, we do the following.

- We remove the subscript j from $\mathbf{D}_{j,\text{new},k}$, $\mathbf{E}_{j,\text{new}}$, and $\zeta_{j,\text{new},k}$ etc. and from everything in Definitions 2.6.3, 2.6.4, 2.6.5, 2.6.9 and 2.7.1.

- Similarly we also let $X_k := X_{\hat{u}_j+k}$ and $\Gamma_k := \Gamma_{j,k}^{\hat{u}_j}$ for both $\hat{u}_j = u_j$ and $\hat{u}_j = u_{j+1}$. More precisely, whenever we say $\mathbb{P}(\text{event} | X_{k-1} \in \Gamma_{k-1}) \geq p_0$ we mean $\mathbb{P}(\text{event} | X_{u_j+k-1} \in \Gamma_{j,k-1}^{u_j}) \geq p_0$ and $\mathbb{P}(\text{event} | X_{u_{j+1}+k-1} \in \Gamma_{j,k-1}^{u_{j+1}}) \geq p_0$.
- Finally, \sum_t refers to $\sum_{t \in \mathcal{J}_u}$ for $u = \hat{u}_j + k$

Also, note the following.

- The proof for the bound on \mathbf{A}_u for $u = u_j + 1$ is the same as that for $u = \hat{u}_j + 1$ since in both cases $\hat{\mathbf{P}}_{t,*} = \hat{\mathbf{P}}_{(j),*}$ and $\hat{\mathbf{P}}_{t,\text{new}} = [\cdot]$ for all $t \in \mathcal{J}_u$. The same is true for the bounds on $\mathbf{A}_{u_j+1,\perp}$ and \mathcal{H}_{u_j+1} .

Fact 2.7.4. When $X_{k-1} \in \Gamma_{k-1}$,

1. $\|\mathbf{D}_{*,k-1}\|_2 \leq \zeta_{j,*}^+$ for $k = 1, \dots, K$.
2. $\|\mathbf{D}_{\text{new},k-1}\|_2 \leq \zeta_{\text{new},k-1}^+$ for $k = 1, \dots, K+1$ (by definition of Γ_{k-1}).
3. Recall that $\zeta_{\text{new},0}^+ = 1$.
4. $\|[(\Phi_t)_{\mathcal{T}_t}'(\Phi_t)_{\mathcal{T}_t}]^{-1}\|_2 \leq \phi^+$ (from Lemma 2.6.15)
5. $\lambda_{\min}(\mathbf{R}_{\text{new}}\mathbf{R}_{\text{new}}') \geq 1 - (\zeta_*^+)^2$ (this follows because $\|\hat{\mathbf{P}}_*'\mathbf{P}_{\text{new}}\|_2 = \|\hat{\mathbf{P}}_*(\mathbf{I} - \mathbf{P}_*\mathbf{P}_*)'\mathbf{P}_{\text{new}}\|_2 \leq \zeta_*$)
6. $\mathbf{E}_{\text{new}}'\mathbf{D}_{\text{new}} = \mathbf{E}_{\text{new}}'\mathbf{E}_{\text{new}}\mathbf{R}_{\text{new}} = \mathbf{R}_{\text{new}}$ and $\mathbf{E}_{\text{new},\perp}'\mathbf{D}_{\text{new}} = \mathbf{0}$.
7. $\tilde{\ell}_t = \mathbf{D}_*\mathbf{a}_{t,*} + \mathbf{D}_{\text{new}}\mathbf{a}_{t,\text{new}}$.
8. \mathbf{e}_t satisfies (2.18) with probability one, i.e. $\mathbf{e}_t = \mathbf{I}_{\mathcal{T}_t}[(\Phi_t)_{\mathcal{T}_t}'(\Phi_t)_{\mathcal{T}_t}]^{-1}\mathbf{I}_{\mathcal{T}_t}'(\mathbf{D}_{*,k-1}\mathbf{a}_{t,*} + \mathbf{D}_{\text{new},k-1}\mathbf{a}_{t,\text{new}})$.

2.7.2 Preliminaries

First observe that the matrices $\mathbf{D}_{\text{new}}, \mathbf{R}_{\text{new}}, \mathbf{E}_{\text{new}}, \mathbf{D}_*, \mathbf{D}_{\text{new},k-1}$ are all functions of the random variable X_{k-1} . Since X_{k-1} is independent of any \mathbf{a}_t for $t \in \mathcal{J}_{\hat{u}_j+k}$, the same is true for these matrices. All terms that we bound for Lemmas 2.6.21 and 2.6.22 are of the form

$\frac{1}{\alpha} \sum_{t \in \mathcal{J}_{\hat{u}_j+k}} \mathbf{Z}_t$ where $\mathbf{Z}_t = f_1(X_{k-1})\mathbf{Y}_t f_2(X_{k-1})$, \mathbf{Y}_t is a sub-matrix of $\mathbf{a}_t \mathbf{a}_t'$, and $f_1(\cdot)$ and $f_2(\cdot)$ are functions of X_{k-1} . Thus, conditioned on X_{k-1} , the \mathbf{Z}_t 's are mutually independent.

All the terms that we bound for Lemma 2.6.23 contain \mathbf{e}_t . Using Lemma 2.6.15, conditioned on X_{k-1} , \mathbf{e}_t satisfies (2.18) with probability one whenever $X_{k-1} \in \Gamma_{k-1}$. Using (2.18), it is easy to see that all the terms needed for this lemma are also of the above form whenever $X_{k-1} \in \Gamma_{k-1}$. Thus, conditioned on X_{k-1} , the \mathbf{Z}_t 's for all the above terms are mutually independent, whenever $X_{k-1} \in \Gamma_{k-1}$.

We will use the following corollaries of the matrix Hoeffding inequality from [23]. These are proved in [11].

Corollary 2.7.5 (Matrix Hoeffding conditioned on another random variable for a nonzero mean Hermitian matrix [23, 11]). *Given an α -length sequence $\{\mathbf{Z}_t\}$ of random Hermitian matrices of size $n \times n$, a r.v. X , and a set \mathcal{C} of values that X can take. Assume that, for all $X \in \mathcal{C}$, (i) \mathbf{Z}_t 's are conditionally independent given X ; (ii) $\mathbb{P}(b_1 \mathbf{I} \preceq \mathbf{Z}_t \preceq b_2 \mathbf{I} | X) = 1$ and (iii) $b_3 \mathbf{I} \preceq \frac{1}{\alpha} \sum_t \mathbb{E}(\mathbf{Z}_t | X) \preceq b_4 \mathbf{I}$. Then for all $\epsilon > 0$,*

$$\mathbb{P} \left(\lambda_{\max} \left(\frac{1}{\alpha} \sum_t \mathbf{Z}_t \right) \leq b_4 + \epsilon \middle| X \right) \geq 1 - n \exp \left(\frac{-\alpha \epsilon^2}{8(b_2 - b_1)^2} \right) \text{ for all } X \in \mathcal{C}$$

$$\mathbb{P} \left(\lambda_{\min} \left(\frac{1}{\alpha} \sum_t \mathbf{Z}_t \right) \geq b_3 - \epsilon \middle| X \right) \geq 1 - n \exp \left(\frac{-\alpha \epsilon^2}{8(b_2 - b_1)^2} \right) \text{ for all } X \in \mathcal{C}$$

Corollary 2.7.6 (Matrix Hoeffding conditioned on another random variable for an arbitrary nonzero mean matrix). *Given an α -length sequence $\{\mathbf{Z}_t\}$ of random matrices of size $n_1 \times n_2$, a r.v. X , and a set \mathcal{C} of values that X can take. Assume that, for all $X \in \mathcal{C}$, (i) \mathbf{Z}_t 's are conditionally independent given X ; (ii) $\mathbb{P}(\|\mathbf{Z}_t\|_2 \leq b_1 | X) = 1$ and (iii) $\|\frac{1}{\alpha} \sum_t \mathbb{E}(\mathbf{Z}_t | X)\|_2 \leq b_2$. Then, for all $\epsilon > 0$,*

$$\mathbb{P} \left(\left\| \frac{1}{\alpha} \sum_t \mathbf{Z}_t \right\|_2 \leq b_2 + \epsilon \middle| X \right) \geq 1 - (n_1 + n_2) \exp \left(\frac{-\alpha \epsilon^2}{32b_1^2} \right) \text{ for all } X \in \mathcal{C}$$

2.7.3 Simple Lemmas Needed for the Proofs

Lemma 2.7.7. *For $j = 1, \dots, J$ and $k = 1, \dots, K$, for all $X_{\hat{u}_j+k-1} \in \Gamma_{j,k-1}^{\hat{u}_j}$*

$$1. \mathbf{0} \preceq \mathbb{E} [\mathbf{a}_{t,*} \mathbf{a}_{t,*}' \mid X_{\hat{u}_j+k-1}] = \mathbf{\Lambda}_{t,*} \preceq \lambda^+ \mathbf{I}$$

$$2. \lambda_{\text{new}}^- \mathbf{I} \preceq \mathbb{E} [\mathbf{a}_{t,\text{new}} \mathbf{a}_{t,\text{new}}' \mid X_{\hat{u}_j+k-1}] = \mathbf{\Lambda}_{t,\text{new}} \preceq \lambda_{\text{new}}^+ \mathbf{I} \text{ and } \hat{\lambda}_{\text{train}}^- \leq \lambda_{\text{new}}^- \leq \lambda_{\text{new}}^+ \leq 3\hat{\lambda}_{\text{train}}^-$$

$$3. \mathbb{E} [\mathbf{a}_{t,*} \mathbf{a}_{t,\text{new}}' \mid X_{\hat{u}_j+k-1}] = \mathbf{0}$$

with $\hat{u}_j = u_j$ or $\hat{u}_j = u_j + 1$.

The same bounds also hold for summation over $t \in \mathcal{J}_{u_j+1}$ when we condition on $X_{u_j} \in \Gamma_{j-1,\text{end}}$.

Proof. The proof follows from Model 2.2.2 and Fact 2.6.13. The only reason we need $X_{\hat{u}_j+k-1} \in \Gamma_{j,k-1}^{\hat{u}_j}$ is to apply Fact 2.6.13 which allows us to lower and upper bound in the eigenvalues of $\mathbf{\Lambda}_{t,\text{new}}$ by λ_{new}^- and λ_{new}^+ and then use (3b). \square

Lemma 2.7.8. Assume that the assumptions of Theorem 2.2.7 hold. Recall that $\mathbf{D}_{\text{new}} = \mathbf{D}_{\text{new},0}$. Conditioned on $X_{k-1} \in \Gamma_{k-1}$,

$$\|\mathbf{I}_{\mathcal{T}}' \mathbf{D}_{\text{new}}\|_2 \leq \kappa_s^+ := .0215 \quad (2.25)$$

for all \mathcal{T} such that $|\mathcal{T}| \leq s$.

The proof is in Appendix 2.B.

2.7.4 Proofs of Lemma 2.6.21 and 2.6.22

Proof of Lemma 2.6.21. We obtain the bounds on \mathbf{A}_u for $u = \hat{u}_j + k$ for $k = 1, 2, \dots, K$ and $\hat{u}_j = u_j$ or $u_j + 1$. For $u = \hat{u}_j + k$, recall that $\mathbf{A}_u := \frac{1}{\alpha} \sum_t \mathbf{E}_{\text{new}}' \tilde{\ell}_t \tilde{\ell}_t' \mathbf{E}_{\text{new}}$.

Notice that $\mathbf{E}_{\text{new}}' \tilde{\ell}_t = \mathbf{R}_{\text{new}} \mathbf{a}_{t,\text{new}} + \mathbf{E}_{\text{new}}' \mathbf{D}_* \mathbf{a}_{t,*}$. Let $\mathbf{Z}_t = \mathbf{R}_{\text{new}} \mathbf{a}_{t,\text{new}} \mathbf{a}_{t,\text{new}}' \mathbf{R}_{\text{new}}'$, and let $\mathbf{Y}_t = \mathbf{R}_{\text{new}} \mathbf{a}_{t,\text{new}} \mathbf{a}_{t,*}' \mathbf{D}_* \mathbf{E}_{\text{new}} + \mathbf{E}_{\text{new}}' \mathbf{D}_* \mathbf{a}_{t,*} \mathbf{a}_{t,\text{new}}' \mathbf{R}_{\text{new}}'$, then

$$\mathbf{A}_u \succeq \frac{1}{\alpha} \sum_t \mathbf{Z}_t + \frac{1}{\alpha} \sum_t \mathbf{Y}_t \quad (2.26)$$

Consider $\frac{1}{\alpha} \sum_t \mathbf{Z}_t$. (1) The \mathbf{Z}_t 's are conditionally independent given X_{k-1} . (2) With probability 1, $\|\mathbf{Z}_t\|_2 \leq r_{\text{new}} \gamma_{\text{new}}^2$. (3) Using a theorem of Ostrowski [24, Theorem 4.5.9], conditioned on $X_{k-1} \in \Gamma_{k-1}$, $\lambda_{\min}(\mathbb{E}[\frac{1}{\alpha} \sum_t \mathbf{Z}_t \mid X_{k-1}]) = \lambda_{\min}(\mathbf{R}_{\text{new}} (\frac{1}{\alpha} \sum_t \mathbf{\Lambda}_{t,\text{new}}) \mathbf{R}_{\text{new}}') \geq \lambda_{\min}(\mathbf{R}_{\text{new}} \mathbf{R}_{\text{new}}') \lambda_{\min}(\frac{1}{\alpha} \sum_t \mathbf{\Lambda}_{t,\text{new}}) \geq (1 - (\zeta_*^+)^2) \lambda_{\text{new}}^-$. The last inequality uses Lemma 2.7.7 and Fact 2.7.4.

Thus, applying Corollary 2.7.5 with ϵ given by (2.20), we get that, for all $X_{k-1} \in \Gamma_{k-1}$,

$$\mathbb{P} \left(\lambda_{\min} \left(\frac{1}{\alpha} \sum_t \mathbf{Z}_t \right) \geq (1 - (\zeta_*^+)^2) \lambda_{\text{new}}^- - \epsilon \middle| X_{k-1} \right) \geq 1 - r_{\text{new}} \exp \left(\frac{-\alpha \zeta^2 (\hat{\lambda}_{\text{train}}^-)^2}{8 \cdot 100^2 \cdot \gamma_{\text{new}}^4} \right). \quad (2.27)$$

Consider $\mathbf{Y}_t = \mathbf{R}_{\text{new}} \mathbf{a}_{t,\text{new}} \mathbf{a}_{t,*}' \mathbf{D}_*' \mathbf{E}_{\text{new}} + \mathbf{E}_{\text{new}}' \mathbf{D}_* \mathbf{a}_{t,*} \mathbf{a}_{t,\text{new}}' \mathbf{R}_{\text{new}}'$. (1) The \mathbf{Y}_t 's are conditionally independent given X_{k-1} . (2) Using the bound on ζ from the theorem, $\|\mathbf{Y}_t\| \leq 2\sqrt{r_{\text{new}} r} \zeta_*^+ \gamma_{\text{new}} \leq 2\sqrt{r_{\text{new}} r} \zeta_*^+ \gamma^2 \leq 2$ holds with probability one for all $X_{k-1} \in \Gamma_{k-1}$. Thus, under the same conditioning, $-2\mathbf{I} \preceq \mathbf{Y}_t \preceq 2\mathbf{I}$ with probability one. (3) By Lemma 2.7.7, $\mathbb{E} \left(\frac{1}{\alpha} \sum_t \mathbf{Y}_t | X_{k-1} \right) = \mathbf{0}$ for all $X_{k-1} \in \Gamma_{k-1}$.

Thus, applying Corollary 2.7.5 with ϵ given by (2.20), we get that, for all $X_{k-1} \in \Gamma_{k-1}$

$$\mathbb{P} \left(\lambda_{\min} \left(\frac{1}{\alpha} \sum_t \mathbf{Y}_t \right) \geq -\epsilon \middle| X_{k-1} \right) \geq 1 - c \exp \left(\frac{-\alpha r_{\text{new}}^2 \zeta^2 (\hat{\lambda}_{\text{train}}^-)^2}{8 \cdot 100^2 \cdot (4)^2} \right) \quad (2.28)$$

Combining (2.26), (2.27) and (2.28) and using the union bound, we get the lemma. \square

Proof of Lemma 2.6.22. Remark 2.7.3 applies.

We obtain the bounds on $\mathbf{A}_{u,\perp}$ for $u = \hat{u}_j + k$ for $k = 1, 2, \dots, K$ with $\hat{u}_j = u_j$ or $u_j + 1$. For all these u 's, recall that $\mathbf{A}_{u,\perp} := \frac{1}{\alpha} \sum_t \mathbf{E}_{\text{new},\perp} \tilde{\ell}_t \tilde{\ell}_t' \mathbf{E}_{\text{new},\perp}$. Using $\mathbf{E}_{\text{new},\perp}' \mathbf{D}_{\text{new}} = 0$, we get that $\mathbf{E}_{\text{new},\perp}' \tilde{\ell}_t = \mathbf{E}_{\text{new},\perp}' \mathbf{D}_* \mathbf{a}_{t,*}$. Thus, $\mathbf{A}_{u,\perp} = \frac{1}{\alpha} \sum_t \mathbf{Z}_t$ with $\mathbf{Z}_t = \mathbf{E}_{\text{new},\perp}' \mathbf{D}_* \mathbf{a}_{t,*} \mathbf{a}_{t,*}' \mathbf{D}_*' \mathbf{E}_{\text{new},\perp}$.

Using the same ideas as for the previous proof we can show that $\mathbf{0} \preceq \mathbf{Z}_t \preceq r(\zeta_*^+)^2 \gamma^2 \mathbf{I} \preceq \zeta \mathbf{I}$ and $\mathbb{E} \left(\frac{1}{\alpha} \sum_t \mathbf{Z}_t | X_{k-1} \right) \preceq (\zeta_*^+)^2 \lambda^+ \mathbf{I}$. Thus by Corollary 2.7.5 the lemma follows. \square

2.7.5 Proof of Lemma 2.6.23

Proof of Lemma 2.6.23. Remark 2.7.3 applies. Using the expression for \mathcal{H}_u given in Definition 2.6.9, and noting that for a basis matrix \mathbf{E} , $\mathbf{E}\mathbf{E}' + \mathbf{E}_\perp \mathbf{E}_\perp' = \mathbf{I}$ we get that

$$\mathcal{H}_u = \frac{1}{\alpha} \sum_{t \in \mathcal{J}_u} \left((\mathbf{I} - \hat{\mathbf{P}}_* \hat{\mathbf{P}}_*') \mathbf{e}_t \mathbf{e}_t' (\mathbf{I} - \hat{\mathbf{P}}_* \hat{\mathbf{P}}_*') - (\tilde{\ell}_t \mathbf{e}_t' (\mathbf{I} - \hat{\mathbf{P}}_* \hat{\mathbf{P}}_*') + (\mathbf{I} - \hat{\mathbf{P}}_* \hat{\mathbf{P}}_*') \mathbf{e}_t \tilde{\ell}_t') + (\mathbf{F}_t + \mathbf{F}_t') \right)$$

where

$$\mathbf{F}_t = \mathbf{E}_{\text{new},\perp} \mathbf{E}_{\text{new},\perp}' \tilde{\ell}_t \tilde{\ell}_t' \mathbf{E}_{\text{new}} \mathbf{E}_{\text{new}}'.$$

Thus,

$$\|\mathcal{H}_u\|_2 \leq 2 \left\| \frac{1}{\alpha} \sum_t \tilde{\ell}_t \mathbf{e}_t' \right\|_2 + \left\| \frac{1}{\alpha} \sum_t \mathbf{e}_t \mathbf{e}_t' \right\|_2 + 2 \left\| \frac{1}{\alpha} \sum_t \mathbf{F}_t \right\|_2 \quad (2.29)$$

Next we obtain high probability bounds on each of the three terms on the right hand side of (2.29).

Consider $\left\| \frac{1}{\alpha} \sum_t \tilde{\ell}_t \mathbf{e}_t' \right\|_2$. Using Lemma 2.6.15, \mathbf{e}_t satisfies (2.18) with probability one for all $X_{k-1} \in \Gamma_{k-1}$.

Let $\mathbf{Z}_t := \tilde{\ell}_t \mathbf{e}_t'$. (1) Conditioned on X_{k-1} , the various \mathbf{Z}_t 's used in the summation are mutually independent, for all $X_{k-1} \in \Gamma_{k-1}$. (2) For $X_{k-1} \in \Gamma_{k-1}$,

$$\|\mathbf{Z}_t\|_2 = \|\tilde{\ell}_t \mathbf{e}_t'\|_2 \leq \left(\zeta_*^+ \sqrt{r} \gamma + \sqrt{r_{\text{new}}} \gamma_{\text{new}} \right) \left(\phi^+ (\zeta_*^+ \sqrt{r} \gamma + \zeta_{\text{new},k-1}^+ \sqrt{r_{\text{new}}} \gamma_{\text{new}}) \right) := b_3$$

holds with probability one. (3) First consider the $k \geq 2$ case. When $X_{k-1} \in \Gamma_{k-1}$,

$$\begin{aligned} & \left\| \mathbb{E} \left[\frac{1}{\alpha} \sum_t \tilde{\ell}_t \mathbf{e}_t' \mid X_{k-1} \right] \right\|_2 \\ &= \left\| \frac{1}{\alpha} \sum_t \left[\left(\mathbf{D}_* \boldsymbol{\Lambda}_{t,*} \mathbf{D}_{*,k-1}' + \mathbf{D}_{\text{new}} \boldsymbol{\Lambda}_{t,\text{new}} \mathbf{D}_{\text{new},k-1}' \right) \mathbf{I}_{\mathcal{T}_t} [(\boldsymbol{\Phi}_t)_{\mathcal{T}_t}' (\boldsymbol{\Phi}_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \right] \right\|_2 \\ &\leq \left[\lambda_{\max} \left(\frac{1}{\alpha} \sum_t \left(\mathbf{D}_* \boldsymbol{\Lambda}_{t,*} \mathbf{D}_{*,k-1}' + \mathbf{D}_{\text{new}} \boldsymbol{\Lambda}_{t,\text{new}} \mathbf{D}_{\text{new},k-1}' \right) \right. \right. \\ &\quad \left. \left. \left(\mathbf{D}_* \boldsymbol{\Lambda}_{t,*} \mathbf{D}_{*,k-1}' + \mathbf{D}_{\text{new}} \boldsymbol{\Lambda}_{t,\text{new}} \mathbf{D}_{\text{new},k-1}' \right)' \right) \right]^{1/2} \\ &\quad \sqrt{\lambda_{\max} \left(\frac{1}{\alpha} \sum_t \left(\mathbf{I}_{\mathcal{T}_t} [(\boldsymbol{\Phi}_t)_{\mathcal{T}_t}' (\boldsymbol{\Phi}_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \right) \left(\mathbf{I}_{\mathcal{T}_t} [(\boldsymbol{\Phi}_t)_{\mathcal{T}_t}' (\boldsymbol{\Phi}_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \right)' \right)} \\ &\leq \left((\zeta_*^+)^2 \lambda^+ + \zeta_{\text{new},k-1}^+ \lambda_{\text{new}}^+ \right) \left(\sqrt{\rho^2 h^+} \phi^+ \right). \end{aligned}$$

The first inequality is by Cauchy-Schwarz for a sum of matrices. This can be found as Lemma 2.D.2 in Appendix 2.D. The second inequality uses Fact 2.7.4 (for the first term of the product) and Lemma 2.5.3 with $\sigma^+ = (\phi^+)^2$ (for the second term of the product).

Now consider the $k = 1$ case. To bound $\left\| \frac{1}{\alpha} \sum_t \mathbf{D}_* \boldsymbol{\Lambda}_{t,*} \mathbf{D}_{*,0}' \mathbf{I}_{\mathcal{T}_t} [(\boldsymbol{\Phi}_t)_{\mathcal{T}_t}' (\boldsymbol{\Phi}_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \right\|_2$ we proceed exactly as we did for the $k \geq 2$ case. We can bound this by $(\zeta_*^+)^2 \lambda^+ \sqrt{\rho^2 h^+} \phi^+$. To bound $\left\| \frac{1}{\alpha} \sum_t \mathbf{D}_{\text{new}} \boldsymbol{\Lambda}_{t,\text{new}} \mathbf{D}_{\text{new},0}' \mathbf{I}_{\mathcal{T}_t} [(\boldsymbol{\Phi}_t)_{\mathcal{T}_t}' (\boldsymbol{\Phi}_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \right\|_2$, we apply Lemma 2.7.8 to get⁷ $\|\mathbf{D}_{\text{new},0}' \mathbf{I}_{\mathcal{T}_t}\|_2 \leq \kappa_s^+$. Using this and Fact 2.7.4, we can bound this by $\kappa_s^+ \lambda_{\text{new}}^+ \phi^+$. Thus, when

⁷Notice that if we want to use the bound of Lemma 2.7.8, we cannot also apply Lemma 2.5.3 for this term. We can get a simpler proof by not using Lemma 2.7.8 at all and proceeding exactly as we did for the $k \geq 2$ case; but doing this will require a much tighter bound on $\rho^2 h^+$ than what we currently need.

$X_0 \in \Gamma_0$,

$$\begin{aligned} & \left\| \mathbb{E} \left[\frac{1}{\alpha} \sum_t \tilde{\ell}_t \mathbf{e}_t' \mid X_0 \right] \right\|_2 \\ &= \left\| \frac{1}{\alpha} \sum_t \left[\left(\mathbf{D}_{*,*} \boldsymbol{\Lambda}_{t,*} \mathbf{D}_{*,*}' + \mathbf{D}_{\text{new}} \boldsymbol{\Lambda}_{t,\text{new}} \mathbf{D}_{\text{new}}' \right) \mathbf{I}_{\mathcal{T}_t} [(\boldsymbol{\Phi}_t)_{\mathcal{T}_t}' (\boldsymbol{\Phi}_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \right] \right\|_2 \\ &\leq \left(\sqrt{\rho^2 h^+} (\zeta_*^+)^2 \lambda^+ + \kappa_s^+ \lambda_{\text{new}}^+ \right) \phi^+. \end{aligned}$$

Thus, by Corollary 2.7.6 with ϵ given by (2.20), we get that, for all $X_{k-1} \in \Gamma_{k-1}$,

$$\mathbb{P} \left(\left\| \frac{1}{\alpha} \sum_t \tilde{\ell}_t \mathbf{e}_t' \right\|_2 \leq b_{\ell e, k} \mid X_{k-1} \right) \geq 1 - n \exp \left(\frac{-\alpha r_{\text{new}}^2 \zeta^2 (\hat{\lambda}_{\text{train}}^-)^2}{32 \cdot 100^2 b_3^2} \right). \quad (2.30)$$

Consider $\|\frac{1}{\alpha} \sum_t \mathbf{e}_t \mathbf{e}_t'\|_2$. Let $\mathbf{Z}_t = \mathbf{e}_t \mathbf{e}_t'$. (1) Conditioned on X_{k-1} , the various \mathbf{Z}_t 's in the summation are independent, for all $X_{k-1} \in \Gamma_{k-1}$. (2) Using Lemma 2.6.15, conditioned on $X_{k-1} \in \Gamma_{k-1}$,

$$\mathbf{0} \preceq \mathbf{Z}_t \preceq \left(\phi^+ (\zeta_*^+ \sqrt{r} \gamma + \zeta_{\text{new}, k-1}^+ \sqrt{r_{\text{new}}} \gamma_{\text{new}}) \right)^2 \mathbf{I} := b_1 \mathbf{I}$$

with probability one. (3) By Fact 2.7.4, when $X_{k-1} \in \Gamma_{k-1}$,

$$\begin{aligned} & \frac{1}{\alpha} \sum_t \mathbb{E} [\mathbf{e}_t \mathbf{e}_t' \mid X_{k-1}] \\ &= \frac{1}{\alpha} \sum_t \mathbf{I}_{\mathcal{T}_t} [(\boldsymbol{\Phi}_t)_{\mathcal{T}_t}' (\boldsymbol{\Phi}_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \left(\mathbf{D}_{*, k-1} \boldsymbol{\Lambda}_{t,*} \mathbf{D}_{*, k-1}' + \right. \\ & \quad \left. \mathbf{D}_{\text{new}, k-1} \boldsymbol{\Lambda}_{t,\text{new}} \mathbf{D}_{\text{new}, k-1}' \right) \mathbf{I}_{\mathcal{T}_t} [(\boldsymbol{\Phi}_t)_{\mathcal{T}_t}' (\boldsymbol{\Phi}_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \end{aligned}$$

When $k = 1$ we can apply Lemma 2.7.8 to get that $\|\mathbf{D}_{\text{new}, 0}' \mathbf{I}_{\mathcal{T}_t}\|_2 \leq \kappa_s^+$. Then we apply Lemma 2.5.3 with $\sigma^+ = (\phi^+)^2 ((\zeta_*^+)^2 \lambda^+ + (\kappa_s^+)^2 \lambda_{\text{new}}^+)$. This gives

$$\mathbf{0} \preceq \mathbb{E} \left[\sum_t \mathbf{e}_t \mathbf{e}_t' \mid X_0 \right] \preceq \rho^2 h^+ (\phi^+)^2 \left((\zeta_*^+)^2 \lambda^+ + (\kappa_s^+)^2 \lambda_{\text{new}}^+ \right) \mathbf{I} \quad \text{for all } X_0 \in \Gamma_0.$$

When $k \geq 2$ we can apply Lemma 2.5.3 with $\sigma^+ = (\phi^+)^2 ((\zeta_*^+)^2 \lambda^+ + (\zeta_{\text{new}, k-1}^+)^2 \lambda_{\text{new}}^+)$ to get that,

$$\mathbf{0} \preceq \mathbb{E} \left[\sum_t \mathbf{e}_t \mathbf{e}_t' \mid X_{k-1} \right] \preceq \rho^2 h^+ (\phi^+)^2 \left((\zeta_*^+)^2 \lambda^+ + (\zeta_{\text{new}, k-1}^+)^2 \lambda_{\text{new}}^+ \right) \mathbf{I} \quad \text{for all } X_{k-1} \in \Gamma_{k-1}.$$

Thus, applying Corollary 2.7.5 with ϵ given by (2.20), we get that, for all $X_{k-1} \in \Gamma_{k-1}$,

$$\mathbb{P} \left(\left\| \frac{1}{\alpha} \sum_t \mathbf{e}_t \mathbf{e}_t' \right\|_2 \leq b_{ee, k} \mid X_{k-1} \right) \geq 1 - n \exp \left(\frac{-\alpha r_{\text{new}}^2 \zeta^2 (\hat{\lambda}_{\text{train}}^-)^2}{8 \cdot 100^2 b_1^2} \right). \quad (2.31)$$

Finally, consider $\left\| \frac{1}{\alpha} \sum_t \mathbf{F}_t \right\|_2$. Since $\mathbf{E}_{\text{new},\perp}' \mathbf{D}_{\text{new}} = 0$,

$$\begin{aligned} \mathbf{F}_t &= \mathbf{E}_{\text{new},\perp} \mathbf{E}_{\text{new},\perp}' \tilde{\ell}_t \tilde{\ell}_t' \mathbf{E}_{\text{new}} \mathbf{E}_{\text{new}}' \\ &= \mathbf{E}_{\text{new},\perp} \mathbf{E}_{\text{new},\perp}' (\mathbf{D}_* \mathbf{a}_{t,*}) (\mathbf{D}_* \mathbf{a}_{t,*} + \mathbf{D}_{\text{new}} \mathbf{a}_{t,\text{new}})' \mathbf{E}_{\text{new}} \mathbf{E}_{\text{new}}' \end{aligned}$$

(1) Conditioned on X_{k-1} , the \mathbf{F}_t 's are mutually independent, for all $X_{k-1} \in \Gamma_{k-1}$. (2) For $X_{k-1} \in \Gamma_{k-1}$,

$$\|\mathbf{F}_t\|_2 \leq (\zeta_*^+)^2 r \gamma^2 + \zeta_*^+ \sqrt{r r_{\text{new}}} \gamma \gamma_{\text{new}} := b_5$$

holds with probability 1. (3) For $X_{k-1} \in \Gamma_{k-1}$,

$$\left\| \mathbb{E} \left[\frac{1}{\alpha} \sum_t \mathbf{F}_t \mid X_{k-1} \right] \right\|_2 \leq \left\| \frac{1}{\alpha} \sum_t (\mathbf{D}_* \mathbf{\Lambda}_{t,*} \mathbf{D}_*') \right\|_2 \leq (\zeta_*^+)^2 \lambda^+ = b_{\mathbf{F}}$$

Applying Corollary 2.7.6 with ϵ given by (2.20), we get that, for all $X_{k-1} \in \Gamma_{k-1}$,

$$\mathbb{P} \left(\left\| \frac{1}{\alpha} \sum_t \mathbf{F}_t \right\|_2 \leq b_{\mathbf{F}} \mid X_{k-1} \right) \geq 1 - n \exp \left(\frac{-\alpha r_{\text{new}}^2 \zeta^2 (\hat{\lambda}_{\text{train}}^-)^2}{32 \cdot 100^2 b_5^2} \right) \quad (2.32)$$

Combining (2.29) with (2.30), (2.31) and (2.32) and using the union bound, we get the lemma. The expression for $p_{\mathcal{H}}$ given in the lemma uses the bounds on ζ from the theorem and uses the loose bound $\zeta_{j,\text{new},k-1}^+ \leq 1$ (to get a simpler expression for the probabilities). \square

2.8 Simulation Experiments

In this section we provide some simulations that demonstrate the robust PCA result we have proven above. More detailed simulations using real data can be found in [17].

The data for Figure 2.7 was generated as follows. We chose $n = 256$ and $t_{\text{max}} = 15,000$. Each measurement had $s = 20$ missing or corrupted entries, i.e. $|\mathcal{T}_t| = 20$. Each non-zero entry of \mathbf{x}_t was drawn uniformly at random between 2 and 6 independent of other entries and other times t . In Figure 2.7 the support of \mathbf{x}_t changes as assumed in Model 2.2.3 with $\rho = 2$ and $\beta = 18$. So the support of \mathbf{x}_t changes by $\frac{s}{2} = 10$ indices every 18 time instants. When the support of \mathbf{x}_t reaches the bottom of the vector, it starts over again at the top. This pattern can be seen in the bottom half of the figure which shows the sparsity pattern of the matrix $\mathbf{S} = [\mathbf{x}_1, \dots, \mathbf{x}_{t_{\text{max}}}]$.

To form the low dimensional vectors ℓ_t , we started with an $n \times r$ matrix of i.i.d. Gaussian entries and orthonormalized the columns. The first $r_0 = 10$ columns of this matrix formed $\mathbf{P}_{(0)}$, the next 2 columns formed $\mathbf{P}_{(1),\text{new}}$, and the last 2 columns formed $\mathbf{P}_{(2),\text{new}}$. We show two subspace changes which occur at $t_1 = 600$ and $t_2 = 8000$. The entries of $\mathbf{a}_{t,*}$ were drawn uniformly at random between -5 and 5, and the entries of $\mathbf{a}_{t,\text{new}}$ were drawn uniformly at random between $-\sqrt{3v_i^{t-t_j}\hat{\lambda}_{\text{train}}^-}$ and $\sqrt{3v_i^{t-t_j}\hat{\lambda}_{\text{train}}^-}$ with $v_i = 1.00017$ and $\hat{\lambda}_{\text{train}}^- = 1$ (and $q_i = 1$). Thus $(\mathbf{A}_{t,\text{new}})_{i,i} = v_i^{t-t_j}\hat{\lambda}_{\text{train}}^-$ as assumed in Model 2.2.2. Entries of \mathbf{a}_t were independent of each other and of the other \mathbf{a}_t 's.

For this simulated data we compare the performance of ReProCS and PCP. The plots show the relative error in recovering ℓ_t , that is $\|\ell_t - \hat{\ell}_t\|_2 / \|\ell_t\|_2$. For the initial subspace estimate $\hat{\mathbf{P}}_0$, we used \mathbf{P}_0 plus some small Gaussian noise and then obtained orthonormal columns. We set $\alpha = 800$ and $K = 6$. For the PCP algorithm, we perform the optimization every α time instants using all of the data up to that point. So the first time PCP is performed on $[\mathbf{m}_1, \dots, \mathbf{m}_\alpha]$ and the second time it is performed on $[\mathbf{m}_1, \dots, \mathbf{m}_{2\alpha}]$ and so on.

Figure 2.7 illustrates the result we have proven. That is ReProCS takes advantage of the initial subspace estimate and slow subspace change (including the bound on γ_{new}) to handle the case when the supports of \mathbf{x}_t are correlated in time. Notice how the ReProCS error increases after a subspace change, but decays exponentially with each projection PCA step. For this data, the PCP program fails to give a meaningful estimate for all but a few times. The average time taken by the ReProCS algorithm was 52 seconds, while PCP averaged over 5 minutes. Simulations were coded in MATLAB[®] and run on a desktop computer with a 3.2 GHz processor.

2.9 Extensions

In this section, we first give other models on changes in \mathcal{T}_t that are special cases of the general model Model 2.5.1 and hence can also be used in Theorem 2.2.5 or 2.2.7. The next three subsections discuss various other results that can also be proved using the proof techniques developed in this work.

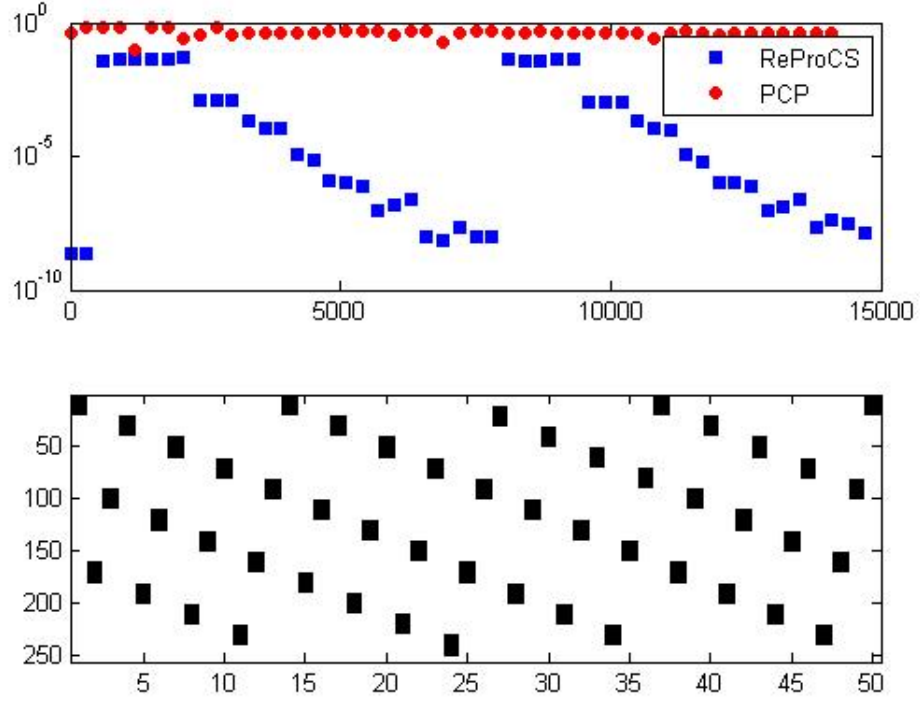


Figure 2.7: Comparison of ReProCS and PCP for the RPCA problem. The top plot is the relative error $\|\ell_t - \hat{\ell}_t\|_2 / \|\ell_t\|_2$. The bottom plot shows the sparsity pattern of \mathcal{S} (black represents a non-zero entry). Results are averaged over 100 simulations and plotted every 300 time instants.

2.9.1 Other Models on Changes in \mathcal{T}_t

We give here other models on changes in \mathcal{T}_t that are special cases of Model 2.5.1.

Model 2.9.1. Suppose that \mathcal{T}_t consists of consecutive indices and is of size s or less, i.e. $|\mathcal{T}_t| \leq s$. When \mathcal{T}_t is not empty, let \tilde{o}_t denote its smallest (topmost) index. Let ρ_1 be an integer. We assume that \tilde{o}_t satisfies the following Bernoulli-Gaussian model:

$$\tilde{o}_t = \lceil o_t \mod n \rceil \text{ where } o_t = o_{t-1} + \theta_t \left(1.1 \frac{s}{\rho} + \varpi_t \right)$$

where $\varpi_t \sim \mathcal{N}(0, \sigma^2)$ (Gaussian) and $\theta_t \sim \text{Bernoulli}(q)$. Assume that $\{\varpi_t\}, \{\theta_t\}$ are mutually independent and independent of ℓ_t 's. Taking the mod with respect to n describes the process of the set \mathcal{T}_t starting over at 1 when its topmost index exceeds n (this models a new object appearing after the old one has disappeared; notice that at any t \mathcal{T}_t could be empty as well, i.e.

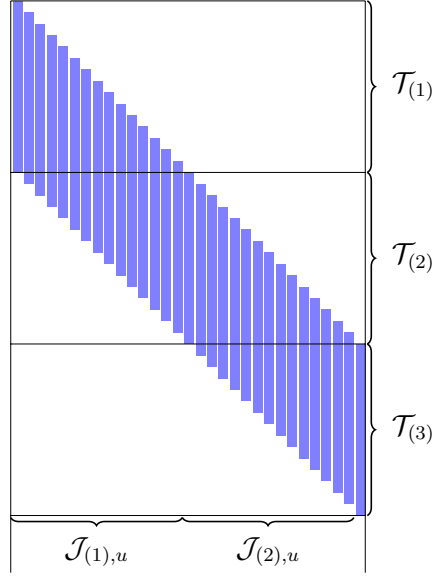


Figure 2.8: Model 2.9.2

there may be no object).

Assume that $s \leq \frac{1.2\rho n}{\alpha}$, $q \geq 1 - \left(\frac{n^{-10}}{2t_{\max}}\right)^{\frac{1}{\beta}}$ for a β that satisfies $\rho^2 \frac{\beta}{\alpha} \leq 0.01$, and $\sigma^2 \leq \frac{s^2}{4000\rho^2 \log(n)}$.

Model 2.9.2. Suppose that \mathcal{T}_t consists of s consecutive indices and suppose that it moves down the vector by between 1 and m indices at every time t . When it reaches the bottom of the vector, we assume that it starts over at 1. Assume that $s \leq 0.0025\alpha$ and $m \leq \frac{n-s}{\alpha}$.

Model 2.9.3. In both models above we let \mathcal{T}_t contain consecutive indices. This models a moving 1D object of length s or less that enters the scene and eventually walks out, and then another object of length s or less may come in. However notice that nothing in our general model, Model 2.5.1, requires the indices to be consecutive or contiguous in any way. Thus in both of Models 2.9.1 and 2.9.2 above, instead of one moving object, we can also have multiple moving objects as long as the union of their supports is of size at most s and satisfies one of these models. Also, with minor changes, the object(s) instead of leaving the scene can reflect back up and start moving in the other direction as well.

Lemma 2.9.4. If $t_{\max} \leq n^{10}$, then Model 2.9.1 is a special case of Model 2.2.3 (and hence a special case of Model 2.5.1) with probability at least $1 - n^{-10}$.

Proof. The proof has three steps. (a) We first use standard arguments about a Bernoulli sequence [25] to prove that the object moves at least once every β time instants with probability at least $1 - 0.5n^{-10}$. The choice of q ensures that this holds. (b) Next we use a standard Gaussian tail bound argument to show that, with probability at least $1 - 0.5n^{-10}$, when it moves, it moves by at least s/ρ indices and at most $1.2s/\rho$ indices. The bound on σ^2 ensures this. (c) The above two claims ensure that, w.h.p., the object remains static for at most β frames at a time and when it moves it moves by at least s/ρ indices and at most $1.2s/\rho$ indices. Notice that all the motion is in one direction. Motion by at least s/ρ in one direction ensures that after the object moves ρ times, i.e. after ρ changes of \mathcal{T}_t , the sets are disjoint, i.e. $\mathcal{T}^{[k]} \cap \mathcal{T}^{[k+\rho]} = \emptyset$. Motion by at most $1.2s/\rho$ in one direction and $1.2\frac{s}{\rho}\alpha \leq n$ ensures the third condition of Model 2.2.3 holds even when the object moves at every frame. \square

Lemma 2.9.5. *Model 2.9.2 is a special case of Model 2.5.1 with $\rho = 2$ and $h^+ = s/\alpha$.*

See Figure 2.8 for a diagram of the model and the idea behind its proof.

Proof. For the sake of clarity, we will prove the case when the object moves exactly 1 index at every time t . The only difference in the general case is the construction of the $\mathcal{J}_{(i),u}$.

Consider an interval \mathcal{J}_u . Let $t_u := (u - 1)\alpha + 1$ denote the first time in \mathcal{J}_u . Without loss of generality (because we can re-label the indices) let the object start at the top of the vector. That is $\mathcal{T}_{t_u} = [1, s]$. Let $l_u = \lceil \frac{n}{s} \rceil$. Let $\mathcal{T}_{(i),u} = [(i - 1)s + 1, is]$ for $i = 1, 2, \dots, \lfloor \frac{n}{s} \rfloor$. If $\frac{n}{s}$ is not an integer, also define $\mathcal{T}_{(\lceil \frac{n}{s} \rceil),u} = [\lfloor \frac{n}{s} \rfloor s + 1, n]$. Define $\mathcal{J}_{(i),u} := [t_u + (i - 1)s, t_u + is - 1]$ for $i = 1, 2, \dots, \lfloor \frac{n}{s} \rfloor$. If $\frac{\alpha}{s}$ is not an integer, also define $\mathcal{J}_{(\lceil \frac{\alpha}{s} \rceil),u} = [t_u + \lfloor \frac{\alpha}{s} \rfloor s, t_u + \alpha - 1]$.

Clearly $\mathcal{J}_{(i),u}$ as defined above are a partition of \mathcal{J}_u . Also, by construction, for all $t \in \mathcal{J}_{(i),u}$, $\mathcal{T}_t \subseteq \mathcal{T}_{(i),u} \cup \mathcal{T}_{(i+1),u}$. This follows from three facts 1) the assumption that $\mathcal{T}_{t_u} = [1, s]$ (which is just a renumbering of the indices to make the numbers clearer) 2) the object moves down by exactly one index at each time t and 3) $m \leq \frac{n-s}{\alpha}$, so that once an index leaves \mathcal{T}_t , it will not return in the next α time instants. A simpler way of stating fact 3) is that the total motion is such that \mathcal{T}_t does not return to where it started i.e. $\mathcal{T}_{t_u} \cap \mathcal{T}_{t_u+\alpha} = \emptyset$.

Notice that $|\mathcal{J}_{(i),u}| \leq s$ for all i . (With the possible exception of the last set, they all have size exactly s .) So under the assumptions of Model 2.9.2 $h_u^*(\alpha) \leq s$, which satisfies Model 2.5.1 with $h^+ = \frac{s}{\alpha} \leq 0.0025\alpha = \frac{0.01\alpha}{2^2} = \frac{0.01\alpha}{\rho^2}$. \square

2.9.2 Analyze the ReProCS algorithm that also removes the deleted directions from the subspace estimate

The tools introduced in this paper – (a) Lemma 2.5.3 and the way it is applied to bound \mathcal{H}_u in Lemma 2.6.23; and (b) the detection lemma (Lemma 2.6.17), the no false detection lemma (Lemma 2.6.16) and the p-PCA lemma (Lemma 2.6.18) – can also be used to get a correctness result for a practical modification of ReProCS with cluster-PCA (ReProCS-cPCA) which is Algorithm 2 of [11]. This algorithm was introduced to also remove the deleted directions from the subspace estimate. It does this by re-estimating the previous subspace at a time after the newly added subspace has been accurately estimated (i.e. at a time after $\hat{t}_j + K\alpha$). A partial result for this algorithm was proved in [11].

This result will need one extra assumption – it will need the eigenvalues of the covariance matrix of ℓ_t to be clustered for a period of time after the subspace change has stabilized, i.e. for a period of d_2 frames in the interval $[t_j + d + 1, t_{j+1} - 1]$ – but it will have a key advantage. It will need a much weaker denseness assumption and hence a much weaker bound on r or r_{mat} . In particular, with this result we expect to be able to allow $r = r_{\text{mat}} \in \mathcal{O}(n)$ with the same assumptions on s and s_{mat} that we currently allow. This requirement is almost as weak as that of PCP.

2.9.3 Relax the independence assumption on ℓ_t 's

The results in this work assume that the ℓ_t 's are independent over time and zero mean; this is a valid model when background images have independent random variations about a fixed mean. Using the tools developed in this paper, a similar result can also be obtained for the more general case of ℓ_t 's following an autoregressive model. This will allow the ℓ_t 's to be correlated over time. A partial result for this case was obtained in [26]. The main change in this case will be that we will need to apply the matrix Azuma inequality from [23] instead

of matrix Hoeffding. This is will also require algebraic manipulation of sums and some other important modifications, as explained in [26], so that the constant term after conditioning on past values of the matrix is small.

2.9.4 Noisy and Undersampled Online Matrix Completion or Online Robust PCA

We expect that the tools introduced in this paper can also be used to analyze the noisy case, i.e. the case of $\mathbf{m}_t = \mathbf{x}_t + \boldsymbol{\ell}_t + \mathbf{w}_t$ where \mathbf{w}_t is small bounded noise. In most practical video applications, while the foreground is truly sparse, the background is only approximately low-rank. The modeling error can be handled as \mathbf{w}_t . The proposed algorithms already apply without modification to this case (see [17] for results on real videos). The reason that our tools will directly extend to the noisy case is this: the sparse recovery step is already a noisy sparse recovery one, its analysis will not change if we also add in more noise due to \mathbf{w}_t . If $\boldsymbol{\ell}_t$ and \mathbf{w}_t are assumed independent, then there should be few simple modifications to the analysis of the p-PCA step as well.

Finally, we expect both the algorithm and the proof techniques to apply with simple changes to the undersampled case $\mathbf{m}_t = \mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \boldsymbol{\ell}_t + \mathbf{w}_t$ as long as \mathbf{B}_t is *not* time-varying, i.e. $\mathbf{B}_t = \mathbf{B}_0$. A partial result for this case was obtained in [27] and experiments were shown in [17].

2.10 Conclusions

In this work, we obtained correctness results for online robust PCA and for online matrix completion. Both results needed four key assumptions: (a) accurate initial subspace knowledge; (b) slow subspace change and mutual independence of the $\boldsymbol{\ell}_t$'s according to Model 2.2.2; (c) *some* changes in the set of missing entries (or in the set of outlier-corrupted entries) over time, one way to quantify what is needed is given in Model 2.2.3; (d) a denseness assumption on the columns of the subspace basis matrices of $\boldsymbol{\ell}_t$; and (e) algorithm parameters are appropriately set.

Ongoing work includes obtaining the results mentioned in Sections 2.9.2, 2.9.3 and 2.9.4. Besides these, we expect the proof techniques developed here to apply to various other problems involving PCA with data and noise terms being correlated.

2.A Appendix A:

Proof that Model 2.2.3 on \mathcal{T}_t satisfies the general Model 2.5.1

Proof of Lemma 2.5.2. Consider an interval \mathcal{J}_u . We will construct one set of mutually disjoint sets $\{\mathcal{T}_{(i),u}\}_{i=1,2,\dots,l_u}$ that are subsets of $\{1, 2, \dots, n\}$ and a partition $\{\mathcal{J}_{(i),u}\}_{i=1,2,\dots,l_u}$ of \mathcal{J}_u so that for all $t \in \mathcal{J}_{(i),u}$, (2.10) holds and so that $h_u(\alpha; \{\mathcal{T}_{(i),u}\}, \{\mathcal{J}_{(i),u}\}) \leq \beta$ for this choice. Since $h_u^*(\alpha)$ takes the minimum over all such sets, this will imply $h_u^*(\alpha) \leq \beta$. By setting $h^+ = \beta/\alpha$ and using the Model 2.2.3 assumption $\rho^2\beta \leq 0.01\alpha$, we will be done.

Recall from Model 2.2.3 that $\mathcal{T}_t = \mathcal{T}^{[k]}$ for all $t \in [t^k, t^{k+1})$ with $t^{k+1} - t^k < \beta$ and $|\mathcal{T}^{[k]}| \leq s$.

Let $t_u := (u-1)\alpha + 1$ denote the first time index of \mathcal{J}_u . Let k_u be the index k for which $t_u \in [t^k, t^{k+1})$. In other words, $\mathcal{T}_{t_u} = \mathcal{T}^{[k_u]}$. Define l_u to be the number of intervals $[t^k, t^{k+1})$ that have non-empty intersection with \mathcal{J}_u . So l_u is one plus the number of times \mathcal{T}_t changes in the interval \mathcal{J}_u . For $i = 1, 2, \dots, l_u - 1$, define

$$\mathcal{T}_{(i),u} := \mathcal{T}^{[k_u+i-1]} \setminus \mathcal{T}^{[k_u+i]},$$

and set $\mathcal{T}_{(l_u),u} = \mathcal{T}^{[k_u+l_u-1]}$. Clearly $l_u \leq \alpha$. Thus, by the Model 2.2.3 assumption (for any k and i such that $k < i \leq k + \alpha$, $(\mathcal{T}^{[k]} \setminus \mathcal{T}^{[k+1]}) \cap (\mathcal{T}^{[i]} \setminus \mathcal{T}^{[i+1]}) = \emptyset$), the $\mathcal{T}_{(i),u}$'s are mutually disjoint.

Next, define a partition of \mathcal{J}_u as

$$\mathcal{J}_{(i),u} := [t^{k_u+i-1}, t^{k_u+i}) \cap \mathcal{J}_u \text{ for } i = 1, 2, \dots, l_u$$

By Model 2.2.3 $1 \leq t_{k+1} - t_k < \beta$ for all k . Since $\mathcal{J}_{(i),u} \subseteq [t^{k_u+i-1}, t^{k_u+i})$, $|\mathcal{J}_{(i),u}| < \beta$ for all $i = 1, 2, \dots, l_u$.

Notice that for all $t \in \mathcal{J}_{(i),u}$, $\mathcal{T}_t = \mathcal{T}^{[k_u+i-1]}$. So if we can show that $\mathcal{T}^{[k_u+i-1]} \subseteq \mathcal{T}_{(i),u} \cup \mathcal{T}_{(i+1),u} \cdots \cup \mathcal{T}_{(i+\rho-1),u}$ for all $i = 1, 2, \dots, l_u$, we will be done since this will imply $h_u^*(\alpha) \leq \beta$.

To show this, set $k = k_u + i - 1$. Then,

$$\begin{aligned}
\mathcal{T}^{[k]} &= \mathcal{T}_{(i),u} \cup [\mathcal{T}^{[k]} \cap \mathcal{T}^{[k+1]}] \\
&= \mathcal{T}_{(i),u} \cup [\mathcal{T}^{[k]} \cap \mathcal{T}^{[k+1]} \setminus \mathcal{T}^{[k+2]}] \cup [\mathcal{T}^{[k]} \cap \mathcal{T}^{[k+1]} \cap \mathcal{T}^{[k+2]}] \\
&\subseteq \mathcal{T}_{(i),u} \cup \mathcal{T}_{(i+1),u} \cup [\mathcal{T}^{[k]} \cap \mathcal{T}^{[k+1]} \cap \mathcal{T}^{[k+2]}] \\
&= \mathcal{T}_{(i),u} \cup \mathcal{T}_{(i+1),u} \cup [\mathcal{T}^{[k]} \cap \mathcal{T}^{[k+1]} \cap \mathcal{T}^{[k+2]} \setminus \mathcal{T}^{[k+3]}] \cup [\mathcal{T}^{[k]} \cap \mathcal{T}^{[k+1]} \cap \mathcal{T}^{[k+2]} \cap \mathcal{T}^{[k+3]}] \\
&\subseteq \mathcal{T}_{(i),u} \cup \mathcal{T}_{(i+1),u} \cup \mathcal{T}_{(i+2),u} \cup [\mathcal{T}^{[k]} \cap \mathcal{T}^{[k+1]} \cap \mathcal{T}^{[k+2]} \cap \mathcal{T}^{[k+3]}].
\end{aligned}$$

Continuing in the same manner as above, we get,

$$\begin{aligned}
\mathcal{T}^{[k]} &\subseteq \mathcal{T}_{(i),u} \cup \mathcal{T}_{(i+1),u} \cup \cdots \cup \mathcal{T}_{(i+\rho-1),u} \cup [\mathcal{T}^{[k]} \cap \mathcal{T}^{[k+1]} \cap \cdots \cap \mathcal{T}^{[k+\rho]}] \\
&= \mathcal{T}_{(i),u} \cup \mathcal{T}_{(i+1),u} \cup \cdots \cup \mathcal{T}_{(i+\rho-1),u}
\end{aligned} \tag{2.33}$$

The last line is because $\mathcal{T}^{[k]} \cap \mathcal{T}^{[k+\rho]} = \emptyset$ by Model 2.2.3. \(\square\)

2.B Appendix B:

Proof of Lemma 2.6.14 (bound on $\zeta_{j,\text{new},k}^+$) and of Lemma 2.7.8

Proof of Lemma 2.6.14. This proof's approach is similar to that of [11, Lemma 6.1]. The details have some differences because our main result now uses different assumptions.

This lemma uses Model 2.5.1. As shown in Lemma 2.5.2, Model 2.2.3 is a special case of this general model.

Recall that $\zeta_{j,\text{new},k}^+ := \frac{b_{\mathcal{H},k}}{b_{\mathcal{A}} - b_{\mathcal{A},\perp} - b_{\mathcal{H},k}}$ with the terms on the RHS defined in Lemmas 2.6.21, 2.6.22, 2.6.23.

Recall that $\epsilon = 0.01r_{\text{new}}\zeta\hat{\lambda}_{\text{train}}^-$. Divide the numerator and denominator by $\hat{\lambda}_{\text{train}}^-$. Define

$$B_k := \begin{cases} \left[\rho^2 h^+ (\phi^+)^2 (\kappa_s^+)^2 (\zeta_{j,\text{new},k-1}^+) + 2\kappa_s^+ \phi^+ \right] \left(\frac{\lambda_{\text{new}}^+}{\hat{\lambda}_{\text{train}}^-} \right) & k = 1 \\ \left[\rho^2 h^+ (\phi^+)^2 \zeta_{j,\text{new},k-1}^+ + 2\sqrt{\rho^2 h^+} \phi^+ \right] \left(\frac{\lambda_{\text{new}}^+}{\hat{\lambda}_{\text{train}}^-} \right) & k \geq 2 \end{cases}$$

$$C_k := \left[\rho^2 h^+ (\phi^+)^2 (\zeta_{j,*}^+) r + 2\sqrt{\rho^2 h^+} \phi^+ (\zeta_{j,*}^+) r + 2(\zeta_{j,*}^+) r \right] \left(\frac{\lambda^+}{\hat{\lambda}_{\text{train}}^-} \right) + 0.05$$

$$D_k := 1 - (\zeta_{j,*}^+)^2 - (\zeta_{j,*}^+)^2 \left(\frac{\lambda^+}{\hat{\lambda}_{\text{train}}^-} \right) - \zeta_{j,\text{new},k-1}^+ B_k - r_{\text{new}} \zeta (C_k + .02)$$

Then,

$$\zeta_{j,\text{new},k}^+ \leq \zeta_{j,\text{new},k-1}^+ \frac{B_k}{D_k} + r_{\text{new}} \zeta \frac{C_k}{D_k}.$$

Recall that $\kappa_s^+ = 0.0215$ and $\phi^+ = 1.2$. It is not difficult to see that $\zeta_{j,\text{new},k}^+$ is an increasing function of $\rho^2 h^+$, r , ζ , $\zeta \frac{\lambda^+}{\hat{\lambda}_{\text{train}}^-}$, and $\frac{\lambda_{\text{new}}^+}{\hat{\lambda}_{\text{train}}^-}$ and $\zeta_{j,\text{new},k-1}^+$. Consider $k = 1$. Using $\zeta_{j,\text{new},0}^+ = 1$ and the upper bounds assumed in Theorem 2.2.7 on the above quantities, we get that $\zeta_{j,\text{new},1}^+ \leq 0.18$.

Thus, $\zeta_{j,\text{new},1}^+ \leq \zeta_{j,\text{new},0}^+ = 1$. Using this and the fact that $\zeta_{j,\text{new},k}^+$ is an increasing function of $\zeta_{j,\text{new},k-1}^+$, we can show by induction that $\zeta_{j,\text{new},k}^+ \leq \zeta_{j,\text{new},k-1}^+$. Thus, $\zeta_{j,\text{new},k}^+ \leq \zeta_{j,\text{new},1}^+ \leq 0.18$ for all $k = 1, 2, \dots, K$.

Using $\zeta_{j,\text{new},k}^+ \leq 0.18$ and the bounds assumed in Theorem 2.2.7 on the other quantities we get that

$$\zeta_{j,\text{new},k}^+ \leq 0.83\zeta_{j,\text{new},k-1}^+ + 0.14r_{\text{new}}\zeta$$

Using this, we get

$$\begin{aligned}
\zeta_{j,\text{new},k}^+ &\leq 0.83\zeta_{j,\text{new},k-1}^+ + 0.14r_{\text{new}}\zeta \leq \zeta_{j,\text{new},0}^+(0.83)^k + \sum_{i=0}^{k-1} (0.83)^i (0.14)r_{\text{new}}\zeta \\
&\leq \zeta_{j,\text{new},0}^+(0.83)^k + \sum_{i=0}^{\infty} (0.83)^i (0.14)r_{\text{new}}\zeta \\
&\leq 0.83^k + 0.84r_{\text{new}}\zeta.
\end{aligned}$$

□

Proof of Lemma 2.7.8. Recall that $\mathbf{D}_{j,\text{new}} = (\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}') \mathbf{P}_{(j),\text{new}}$. Then $\|\mathbf{I}_{\mathcal{T}'} \mathbf{D}_{j,\text{new}}\|_2 = \|\mathbf{I}_{\mathcal{T}'} (\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}') \mathbf{P}_{(j),\text{new}}\|_2 \leq \|\mathbf{I}_{\mathcal{T}'} \mathbf{P}_{(j),\text{new}}\|_2 + \|\hat{\mathbf{P}}_{(j),*}' \mathbf{P}_{(j),\text{new}}\|_2 \leq \kappa_s(\mathbf{P}_{(j),\text{new}}) + \|\hat{\mathbf{P}}_{(j),*}' (\mathbf{I} - \mathbf{P}_{(j),*} \mathbf{P}_{(j),*}') \mathbf{P}_{(j),\text{new}}\|_2 \leq \kappa_s(\mathbf{P}_{(j),\text{new}}) + \zeta_{j,*}$. The event $X_{\hat{u}_j+k-1} \in \Gamma_{j,k-1}^{\hat{u}_j}$ implies that $\zeta_{j,*} \leq \zeta_{j,*}^+ \leq 0.0015$. Thus, the lemma follows. □

2.C Appendix C:

Proof of the Compressed Sensing (CS) Lemma (Lemma 2.6.15)

This proof's approach is similar to that of [11, Lemma 6.4]. The details have some differences because our main result now uses different assumptions. The proof uses the denseness assumption and subspace error bounds $\zeta_{j,*} \leq \zeta_{j,*}^+$ and $\zeta_{j,\text{new},k-1} \leq \zeta_{j,\text{new},k-1}^+$, that hold when $X_{\hat{u}_j+k-1} \in \Gamma_{j,k-1}^{\hat{u}_j}$ for $\hat{u}_j = u_j$ or $\hat{u}_j = u_j + 1$, to obtain bounds on the restricted isometry constant (RIC) of the sparse recovery matrix Φ_t and the sparse recovery error $\|\mathbf{b}_t\|_2$. Applying the noisy compressed sensing (CS) result from [19] and the assumed bounds on ζ and γ , the lemma follows.

Lemma 2.C.1. [11, Lemma 2.10] Suppose that \mathbf{P} , $\hat{\mathbf{P}}$ and \mathbf{Q} are three basis matrices. Also, \mathbf{P} and $\hat{\mathbf{P}}$ are of the same size, $\mathbf{Q}'\mathbf{P} = \mathbf{0}$ and $\|(\mathbf{I} - \hat{\mathbf{P}}\hat{\mathbf{P}}')\mathbf{P}\|_2 = \zeta_*$. Then,

1. $\|(\mathbf{I} - \hat{\mathbf{P}}\hat{\mathbf{P}}')\mathbf{P}\mathbf{P}'\|_2 = \|(\mathbf{I} - \mathbf{P}\mathbf{P}')\hat{\mathbf{P}}\hat{\mathbf{P}}'\|_2 = \|(\mathbf{I} - \mathbf{P}\mathbf{P}')\hat{\mathbf{P}}\|_2 = \|(\mathbf{I} - \hat{\mathbf{P}}\hat{\mathbf{P}}')\mathbf{P}\|_2 = \zeta_*$
2. $\|\mathbf{P}\mathbf{P}' - \hat{\mathbf{P}}\hat{\mathbf{P}}'\|_2 \leq 2\|(\mathbf{I} - \hat{\mathbf{P}}\hat{\mathbf{P}}')\mathbf{P}\|_2 = 2\zeta_*$
3. $\|\hat{\mathbf{P}}'\mathbf{Q}\|_2 \leq \zeta_*$
4. $\sqrt{1 - \zeta_*^2} \leq \sigma_i\left((\mathbf{I} - \hat{\mathbf{P}}\hat{\mathbf{P}}')\mathbf{Q}\right) \leq 1$

We begin by first bounding the RIC of the CS matrix Φ_t . We will use the notation $\kappa_s^2(\mathbf{P})$ to mean $(\kappa_s(\mathbf{P}))^2$.

Lemma 2.C.2 (Bounding the RIC of Φ_t [11, Lemma 6.6]). Recall that $\zeta_{j,*} := \|(\mathbf{I} - \hat{\mathbf{P}}_{(j),*}\hat{\mathbf{P}}_{(j),*}')\mathbf{P}_{(j),*}\|_2$. The following hold.

1. Suppose that a basis matrix \mathbf{P} can be split as $\mathbf{P} = [\mathbf{P}_1 \ \mathbf{P}_2]$ where \mathbf{P}_1 and \mathbf{P}_2 are also basis matrices. Then $\kappa_s^2(\mathbf{P}) = \max_{\mathcal{T}: |\mathcal{T}| \leq s} \|\mathbf{I}_{\mathcal{T}'}\mathbf{P}\|_2^2 \leq \kappa_s^2(\mathbf{P}_1) + \kappa_s^2(\mathbf{P}_2)$.
2. $\kappa_s^2(\hat{\mathbf{P}}_{(j),*}) \leq (\kappa_{s,*})^2 + 2\zeta_*$ for all j
3. $\kappa_s(\hat{\mathbf{P}}_{(j),\text{new},k}) \leq \kappa_{s,\text{new}} + \zeta_{j,\text{new},k} + \zeta_{j,*}$ for all j and k .
4. For $t \in [(u_{j-1} + K)\alpha + 1, (\hat{u}_j + 1)\alpha)$, $\delta_s(\Phi_t) = \kappa_s^2(\hat{\mathbf{P}}_{(j),*}) \leq (\kappa_{s,*})^2 + 2\zeta_{j,*}$.

5. For $k = 1, \dots, K-1$, for $t \in [(\hat{u}_j + k)\alpha + 1, (\hat{u}_j + k + 1)\alpha]$ $\delta_s(\Phi_t) = \kappa_s^2([\hat{P}_{(j),*}, \hat{P}_{(j),\text{new},k}]) \leq \kappa_s^2(\hat{P}_{(j),*}) + \kappa_s^2(\hat{P}_{(j),\text{new},k}) \leq (\kappa_{s,*})^2 + 2\zeta_{j,*} + (\kappa_{s,\text{new}} + \zeta_{j,\text{new},k} + \zeta_{j,*})^2$.

Proof. 1. Recall that $\kappa_s^2(P) = \max_{|\mathcal{T}| \leq s} \|I_{\mathcal{T}}' P\|_2^2$. Also, $\|I_{\mathcal{T}}' P\|_2^2 = \|I_{\mathcal{T}}'[P_1 \ P_2][P_1 \ P_2]' I_{\mathcal{T}}\|_2 = \|I_{\mathcal{T}}'(P_1 P_1' + P_2 P_2') I_{\mathcal{T}}\|_2 \leq \|I_{\mathcal{T}}' P_1 P_1' I_{\mathcal{T}}\|_2 + \|I_{\mathcal{T}}' P_2 P_2' I_{\mathcal{T}}\|_2$. Thus, the inequality follows.

2. For any set \mathcal{T} with $|\mathcal{T}| \leq s$, $\|I_{\mathcal{T}}' \hat{P}_{(j),*}\|_2^2 = \|I_{\mathcal{T}}' \hat{P}_{(j),*} \hat{P}_{(j),*}' I_{\mathcal{T}}\|_2 = \|I_{\mathcal{T}}'(\hat{P}_{(j),*} \hat{P}_{(j),*}' - P_{(j),*} P_{(j),*}' + P_{(j),*} P_{(j),*}') I_{\mathcal{T}}\|_2 \leq \|I_{\mathcal{T}}'(\hat{P}_{(j),*} \hat{P}_{(j),*}' - P_{(j),*} P_{(j),*}') I_{\mathcal{T}}\|_2 + \|I_{\mathcal{T}}' P_{(j),*} P_{(j),*}' I_{\mathcal{T}}\|_2 \leq 2\zeta_{j,*} + (\kappa_{s,*})^2$. The last inequality follows using Lemma 2.C.1 with $P = P_{(j),*}$ and $\hat{P} = \hat{P}_{(j),*}$.

3. By Lemma 2.C.1 with $P = P_{(j),*}$, $\hat{P} = \hat{P}_{(j),*}$ and $Q = P_{(j),\text{new}}$, $\|P_{(j),\text{new}}' \hat{P}_{(j),*}\|_2 \leq \zeta_{j,*}$. By Lemma 2.C.1 with $P = P_{(j),\text{new}}$ and $\hat{P} = \hat{P}_{(j),\text{new},k}$, $\|(I - P_{(j),\text{new}} P_{(j),\text{new}}') \hat{P}_{(j),\text{new},k}\|_2 = \|(I - \hat{P}_{(j),\text{new},k} \hat{P}_{(j),\text{new},k}') P_{(j),\text{new}}\|_2$.

For any set \mathcal{T} with $|\mathcal{T}| \leq s$, $\|I_{\mathcal{T}}' \hat{P}_{(j),\text{new},k}\|_2 \leq \|I_{\mathcal{T}}'(I - P_{(j),\text{new}} P_{(j),\text{new}}') \hat{P}_{(j),\text{new},k}\|_2 + \|I_{\mathcal{T}}' P_{(j),\text{new}} P_{(j),\text{new}}' \hat{P}_{(j),\text{new},k}\|_2 \leq \|(I - P_{(j),\text{new}} P_{(j),\text{new}}') \hat{P}_{(j),\text{new},k}\|_2 + \|I_{\mathcal{T}}' P_{(j),\text{new}}\|_2 = \|(I - \hat{P}_{(j),\text{new},k} \hat{P}_{(j),\text{new},k}') P_{(j),\text{new}}\|_2 + \|I_{\mathcal{T}}' P_{(j),\text{new}}\|_2 \leq \|D_{(j),\text{new},k}\|_2 + \|\hat{P}_{(j),*} \hat{P}_{(j),*}' P_{(j),\text{new}}\|_2 + \|I_{\mathcal{T}}' P_{(j),\text{new}}\|_2$. Taking max over $|\mathcal{T}| \leq s$ the claim follows.

4. This follows using Lemma 2.2.9 and the second claim of this lemma.

5. This follows using Lemma 2.2.9 and the first three claims of this lemma.

□

Corollary 2.C.3.

1. Conditioned on $\Gamma_{j-1,\text{end}}$, for $t \in [t_j, (\hat{u}_j + 1)\alpha]$, $\delta_s(\Phi_t) \leq \delta_{2s}(\Phi_t) \leq (\kappa_{2s,*})^2 + 2\zeta_{j,*}^+ < 0.1 < 0.1479$, and $\|[(\Phi_t)_{\mathcal{T}_t}'(\Phi_t)_{\mathcal{T}_t}]^{-1}\|_2 \leq \frac{1}{1-\delta_s(\Phi_t)} < 1.2 := \phi^+$.
2. For $k = 2, \dots, K$ and $\hat{u}_j = u_j$ or $\hat{u}_j = u_j + 1$, conditioned on $\Gamma_{j,k-1}^{\hat{u}_j}$, for $t \in [(\hat{u}_j + k - 1)\alpha + 1, (\hat{u}_j + k)\alpha]$, $\delta_s(\Phi_t) \leq \delta_{2s}(\Phi_t) \leq (\kappa_{2s,*})^2 + 2\zeta_{j,*}^+ + (\kappa_{2s,\text{new}} + \zeta_{j,\text{new},k-1}^+ + \zeta_{j,*}^+)^2 < 0.1479$, and $\|[(\Phi_t)_{\mathcal{T}_t}'(\Phi_t)_{\mathcal{T}_t}]^{-1}\|_2 \leq \frac{1}{1-\delta_s(\Phi_t)} < 1.2 := \phi^+$.

3. For $\hat{u}_j = u_j$ or $\hat{u}_j = u_j + 1$, conditioned on $\Gamma_{j,K}^{\hat{u}_j}$, for $t \in [(\hat{u}_j + K)\alpha + 1, t_{j+1} - 1]$,
 $\delta_s(\Phi_t) \leq \delta_{2s}(\Phi_t) \leq (\kappa_{2s,*})^2 + 2\zeta_{j,*}^+ < 0.1 < 0.1479$, and $\|[(\Phi_t)_{\mathcal{T}_t}'(\Phi_t)_{\mathcal{T}_t}]^{-1}\|_2 \leq \frac{1}{1-\delta_s(\Phi_t)} < 1.2 := \phi^+$.

Proof. This follows using Lemma 2.C.2, the definitions of $\Gamma_{j-1,\text{end}}$ and $\Gamma_{j,k}^{\hat{u}_j}$, and the bound on $\zeta_{j,\text{new},k-1}^+$ from Lemma 2.6.14. \square

The following are straightforward bounds that will be useful for the proof of Lemma 2.6.15.

Fact 2.C.4. *Under the assumptions of Theorem 2.2.7:*

- $\zeta_{j,*}^+ \gamma \leq \frac{\sqrt{\zeta}}{\sqrt{r_0 + (J-1)c}} \leq \sqrt{\zeta}$
- $\zeta_{j,\text{new},k-1}^+ \leq 0.83^{k-1} + 0.84r_{\text{new}}\zeta$ (from Lemma 2.6.14)
- $\zeta_{j,\text{new},k-1}^+ \gamma_{\text{new}} \leq 0.83^{k-1}\gamma_{\text{new}} + 0.84r_{\text{new}}\zeta\gamma_{\text{new}} \leq 0.83^{k-1}\gamma_{\text{new}} + 0.3\sqrt{\zeta}$

Proof of Lemma 2.6.15. We will prove claim 2). The others are done in the same way.

Recall that $\Gamma_{j,k-1}^{\hat{u}_j}$ implies that $\zeta_{j,*} \leq \zeta_{j,*}^+$ and $\zeta_{j,\text{new},k-1} \leq \zeta_{j,\text{new},k-1}^+$.

- a) For $t \in [(\hat{u}_j + k - 1)\alpha + 1, (\hat{u}_j + k)\alpha]$, $\mathbf{b}_t := (\mathbf{I} - \hat{\mathbf{P}}_{t-1}\hat{\mathbf{P}}_{t-1}')\ell_t = \mathbf{D}_{j,*,k-1}\mathbf{a}_{t,*} + \mathbf{D}_{j,\text{new},k-1}\mathbf{a}_{t,\text{new}}$.

Thus, using Fact 2.C.4

$$\begin{aligned} \|\mathbf{b}_t\|_2 &\leq \zeta_{j,*}\sqrt{r}\gamma + \zeta_{j,\text{new},k-1}\sqrt{r_{\text{new}}}\gamma_{\text{new}} \\ &\leq \sqrt{\zeta}\sqrt{r} + (0.83^{k-1}\gamma_{\text{new}} + 0.84\sqrt{\zeta})\sqrt{r_{\text{new}}} \\ &= \sqrt{r_{\text{new}}}0.83^{k-1}\gamma_{\text{new}} + \sqrt{\zeta}(\sqrt{r} + 0.84\sqrt{r_{\text{new}}}) \leq \xi. \end{aligned}$$

- b) By Corollary 2.C.3, $\delta_{2s}(\Phi_t) < 0.15 < \sqrt{2} - 1$. Given $|\mathcal{T}_t| \leq s$, $\|\mathbf{b}_t\|_2 \leq \xi$, by the theorem in [19], the CS error satisfies

$$\|\hat{\mathbf{x}}_{t,\text{cs}} - \mathbf{x}_t\|_2 \leq \frac{4\sqrt{1 + \delta_{2s}(\Phi_t)}}{1 - (\sqrt{2} + 1)\delta_{2s}(\Phi_t)}\xi < 7\xi.$$

- c) Using the above, $\|\hat{\mathbf{x}}_{t,\text{cs}} - \mathbf{x}_t\|_\infty \leq 7\xi$. Since $\min_{i \in \mathcal{T}_t} |(\mathbf{x}_t)_i| \geq x_{\min}$ and $(\mathbf{x}_t)_{\mathcal{T}_t^c} = 0$, $\min_{i \in \mathcal{T}_t} |(\hat{\mathbf{x}}_{t,\text{cs}})_i| \geq x_{\min} - 7\xi$ and $\max_{i \in \bar{\mathcal{T}}_t} |(\hat{\mathbf{x}}_{t,\text{cs}})_i| \leq 7\xi$. If $\omega < x_{\min} - 7\xi$, then $\hat{\mathcal{T}}_t \supseteq \mathcal{T}_t$. On the other hand, if $\omega > 7\xi$, then $\hat{\mathcal{T}}_t \subseteq \mathcal{T}_t$. Since ω satisfies $7\xi \leq \omega \leq x_{\min} - 7\xi$, the support of \mathbf{x}_t is exactly recovered, i.e. $\hat{\mathcal{T}}_t = \mathcal{T}_t$.

- d) Given $\hat{\mathcal{T}}_t = \mathcal{T}_t$, the least squares estimate of \mathbf{x}_t satisfies $(\hat{\mathbf{x}}_t)_{\mathcal{T}_t} = [(\mathbf{\Phi}_t)_{\mathcal{T}_t}]^\dagger \mathbf{y}_t = [(\mathbf{\Phi}_t)_{\mathcal{T}_t}]^\dagger (\mathbf{\Phi}_t \mathbf{x}_t + \mathbf{\Phi}_t \boldsymbol{\ell}_t)$ and $(\hat{\mathbf{x}}_t)_{\bar{\mathcal{T}}_t} = \mathbf{0}$. Also, $(\mathbf{\Phi}_t)_{\mathcal{T}_t}' \mathbf{\Phi}_t = \mathbf{I}_{\mathcal{T}_t}' \mathbf{\Phi}_t$ (this follows since $(\mathbf{\Phi}_t)_{\mathcal{T}_t} = \mathbf{\Phi}_t \mathbf{I}_{\mathcal{T}_t}$ and $\mathbf{\Phi}_t' \mathbf{\Phi}_t = \mathbf{\Phi}_t$). Using this, the error $\mathbf{e}_t := \hat{\mathbf{x}}_t - \mathbf{x}_t$ satisfies (2.18). Thus, using Fact 2.C.4 and the bounds on $\|\mathbf{a}_t\|_\infty$ and $\|\mathbf{a}_{t,\text{new}}\|_\infty$, for $t \in [(\hat{u}_j + k - 1)\alpha + 1, (\hat{u}_j + k)\alpha]$,

$$\|\mathbf{e}_t\|_2 \leq \phi^+(\zeta_{j,*}^+ \sqrt{r} \gamma + \zeta_{j,\text{new},k-1}^+ \sqrt{r_{\text{new}}} \gamma_{\text{new}}) \leq 1.2 \left(1.06 \sqrt{\zeta} + (0.83)^{k-1} \sqrt{r_{\text{new}}} \gamma_{\text{new}} \right)$$

The last inequality follows from Lemma 2.6.14.

□

2.D Appendix D: Proof of Cauchy-Schwarz inequality for matrices

Lemma 2.D.1 (Cauchy-Schwarz for a sum of vectors). *For vectors \mathbf{x}_t and \mathbf{y}_t ,*

$$\left(\sum_{t=1}^{\alpha} \mathbf{x}_t' \mathbf{y}_t \right)^2 \leq \left(\sum_t \|\mathbf{x}_t\|_2^2 \right) \left(\sum_t \|\mathbf{y}_t\|_2^2 \right)$$

Proof.

$$\begin{aligned} \left(\sum_{t=1}^{\alpha} \mathbf{x}_t' \mathbf{y}_t \right)^2 &= \left([\mathbf{x}_1', \dots, \mathbf{x}_{\alpha}'] \begin{bmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_{\alpha} \end{bmatrix} \right)^2 \leq \left\| \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_{\alpha} \end{bmatrix} \right\|_2^2 \left\| \begin{bmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_{\alpha} \end{bmatrix} \right\|_2^2 \\ &= \left(\sum_{t=1}^{\alpha} \|\mathbf{x}_t\|_2^2 \right) \left(\sum_{t=1}^{\alpha} \|\mathbf{y}_t\|_2^2 \right) \end{aligned}$$

The inequality is by Cauchy-Schwarz for a single vector. \square

Lemma 2.D.2 (Cauchy-Schwarz for a sum of matrices). *For matrices \mathbf{X}_t and \mathbf{Y}_t ,*

$$\left\| \frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbf{X}_t \mathbf{Y}_t' \right\|_2^2 \leq \lambda_{\max} \left(\frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbf{X}_t \mathbf{X}_t' \right) \lambda_{\max} \left(\frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbf{Y}_t \mathbf{Y}_t' \right)$$

Proof of Lemma 2.D.2.

$$\begin{aligned} \left\| \sum_{t=1}^{\alpha} \mathbf{X}_t \mathbf{Y}_t' \right\|_2^2 &= \max_{\substack{\|\mathbf{x}\|=1 \\ \|\mathbf{y}\|=1}} \left| \mathbf{x}' \left(\sum_t \mathbf{X}_t \mathbf{Y}_t' \right) \mathbf{y} \right|^2 \\ &= \max_{\substack{\|\mathbf{x}\|=1 \\ \|\mathbf{y}\|=1}} \left| \sum_{t=1}^{\alpha} (\mathbf{X}_t' \mathbf{x})' (\mathbf{Y}_t' \mathbf{y}) \right|^2 \\ &\leq \max_{\substack{\|\mathbf{x}\|=1 \\ \|\mathbf{y}\|=1}} \left(\sum_{t=1}^{\alpha} \|\mathbf{X}_t' \mathbf{x}\|_2^2 \right) \left(\sum_{t=1}^{\alpha} \|\mathbf{Y}_t' \mathbf{y}\|_2^2 \right) \\ &= \max_{\|\mathbf{x}\|=1} \mathbf{x}' \sum_{t=1}^{\alpha} \mathbf{X}_t \mathbf{X}_t' \mathbf{x} \cdot \max_{\|\mathbf{y}\|=1} \mathbf{y}' \sum_{t=1}^{\alpha} \mathbf{Y}_t \mathbf{Y}_t' \mathbf{y} \\ &= \lambda_{\max} \left(\sum_{t=1}^{\alpha} \mathbf{X}_t \mathbf{X}_t' \right) \lambda_{\max} \left(\sum_{t=1}^{\alpha} \mathbf{Y}_t \mathbf{Y}_t' \right) \end{aligned}$$

The inequality is by Lemma 2.D.1. The penultimate line is because $\|\mathbf{x}\|_2^2 = \mathbf{x}' \mathbf{x}$. Multiplying both sides by $(\frac{1}{\alpha})^2$ gives the desired result. \square

References

- [1] E. J. Candès, X. Li, Y. Ma, and J. Wright, “Robust principal component analysis?” *Journal of ACM*, vol. 58, no. 3, 2011.
- [2] V. Chandrasekaran, S. Sanghavi, P. A. Parrilo, and A. S. Willsky, “Rank-sparsity incoherence for matrix decomposition,” *SIAM Journal on Optimization*, vol. 21, 2011.
- [3] D. Hsu, S. Kakade, and T. Zhang, “Robust matrix decomposition with sparse corruptions,” *IEEE Trans. Info. Th.*, Nov. 2011.
- [4] M. Fazel, “Matrix rank minimization with applications,” *PhD thesis, Stanford University*, 2002.
- [5] E. J. Candès and B. Recht, “Exact matrix completion via convex optimization,” *Found. of Comput. Math.*, no. 9, pp. 717–772, 2008.
- [6] M. Brand, “Incremental singular value decomposition of uncertain data with missing values,” in *Eur. Conf. on Comp. Vis. (ECCV)*, 2002.
- [7] L. Balzano, B. Recht, and R. Nowak, “Online Identification and Tracking of Subspaces from Highly Incomplete Information,” in *Allerton Conf. Communication, Control, and Computing*, 2010.
- [8] Y. Chi, Y. C. Eldar, and R. Calderbank, “Petrels: Parallel subspace estimation and tracking by recursive least squares from partial observations,” *IEEE Trans. Sig. Proc.*, December 2013.
- [9] L. Balzano and S. Wright, “Local convergence of an algorithm for subspace identification from partial data,” *arXiv:1306.3391 [cs.NA]*.
- [10] A. Krishnamurthy and A. Singh, “Low-rank matrix and tensor completion via adaptive sampling,” in *Advances in Neural Information Processing Systems 26*, C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Weinberger, Eds. Curran Associates, Inc., 2013, pp. 836–844. [Online]. Available: <http://papers.nips.cc/paper/4954-low-rank-matrix-and-tensor-completion-via-adaptive-sampling.pdf>
- [11] C. Qiu, N. Vaswani, B. Loos, and L. Hogben, “Recursive robust pca or recursive sparse recovery in large but structured noise,” *IEEE Trans. Info. Th.*, Aug. 2014, shorter versions in ICASSP 2013 and ISIT 2013.
- [12] J. Feng, H. Xu, and S. Yan, “Online robust pca via stochastic optimization,” in *Adv. Neural Info. Proc. Sys. (NIPS)*, 2013.

- [13] J. Feng, H. Xu, S. Mannor, and S. Yan, "Online pca for contaminated data," in *Advances in Neural Information Processing Systems 26*, C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Weinberger, Eds. Curran Associates, Inc., 2013, pp. 764–772. [Online]. Available: <http://papers.nips.cc/paper/5135-online-pca-for-contaminated-data.pdf>
- [14] J. He, L. Balzano, and A. Szlam, "Incremental gradient on the grassmannian for online foreground and background separation in subsampled video," in *IEEE Conf. on Comp. Vis. Pat. Rec. (CVPR)*, 2012.
- [15] P. Netrapalli, P. Jain, and S. Sanghavi, "Low-rank matrix completion using alternating minimization," in *Symposium on Theory of Computing (STOC)*, 2013.
- [16] B. Lois and N. Vaswani, "A correctness result for online robust pca," in *IEEE Intl. Conf. Acoustics, Speech, Sig. Proc. (ICASSP)*, 2015.
- [17] H. Guo, C. Qiu, and N. Vaswani, "An online algorithm for separating sparse and low-dimensional signal sequences from their sum," *IEEE Trans. Sig. Proc.*, pp. 4284–4297, Aug. 2014.
- [18] F. D. L. Torre and M. J. Black, "A framework for robust subspace learning," *International Journal of Computer Vision*, vol. 54, pp. 117–142, 2003.
- [19] E. Candes, "The restricted isometry property and its implications for compressed sensing," *Compte Rendus de l'Academie des Sciences, Paris, Serie I*, pp. 589–592, 2008.
- [20] J. Feng, H. Xu, and S. Yan, "Online robust pca via stochastic optimization," in *Advances in Neural Information Processing Systems 26*, C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Weinberger, Eds. Curran Associates, Inc., 2013, pp. 404–412. [Online]. Available: <http://papers.nips.cc/paper/5131-online-robust-pca-via-stochastic-optimization.pdf>
- [21] J. Zhan and N. Vaswani, "Robust pca with partial subspace knowledge," in *IEEE Intl. Symp. Info. Th. (ISIT)*, 2014.
- [22] C. Davis and W. M. Kahan, "The rotation of eigenvectors by a perturbation. iii," *SIAM Journal on Numerical Analysis*, Mar. 1970.
- [23] J. A. Tropp, "User-friendly tail bounds for sums of random matrices," *Foundations of Computational Mathematics*, vol. 12, no. 4, 2012.
- [24] R. Horn and C. Johnson, *Matrix Analysis*. Cambridge Univ. Press, 1985.
- [25] M. Muselli, "On convergence properties of pocket algorithm," *Neural Networks, IEEE Transactions on*, vol. 8, no. 3, pp. 623–629, May 1997.
- [26] J. Zhan and N. Vaswani, "Performance guarantees for reprocs – correlated low-rank matrix entries case," arXiv:1405.5887 [cs.IT].
- [27] B. Lois, N. Vaswani, and C. Qiu, "Performance guarantees for undersampled recursive sparse recovery in large but structured noise," in *GlobalSIP*, 2013.

CHAPTER 3. RECURSIVE SPARSE RECOVERY IN CORRELATED STRUCTURED NOISE

A paper prepared for submission to *IEEE Transactions on Information Theory*

Brian Lois, Namrata Vaswani, and Jinchun Zhan

Abstract

This work studies the problem of recursive robust principal components analysis (PCA). At each time t , suppose that a vector $\mathbf{m}_t = \boldsymbol{\ell}_t + \mathbf{x}_t$ is observed. The vectors $\boldsymbol{\ell}_t$ lie in a slowly changing and dense low-dimensional subspace. The vectors \mathbf{x}_t are sparse and their support changes frequently enough. The goal is to recover \mathbf{x}_t and $\boldsymbol{\ell}_t$ at each time t and to maintain an estimate of $\text{range}(\boldsymbol{\ell}_1, \dots, \boldsymbol{\ell}_t)$. This work improves on existing results by assuming a more general autoregressive model for the low-dimensional vectors $\boldsymbol{\ell}_t$. It is proven that under certain model assumptions, with high probability, the practical ReProCS algorithm exactly recovers the support of \mathbf{x}_t , and the error made in estimating \mathbf{x}_t and $\boldsymbol{\ell}_t$ is bounded by a small value that depends on the accuracy of the initial subspace estimate. Also, all of the subspace changes are detected within a certain delay, and the error made in estimating the new subspace decays below a small value within a finite delay.

3.1 Introduction

Principal Components Analysis (PCA) is a tool that is frequently used for dimension reduction. Given a matrix of data \mathbf{D} , PCA seeks to recover a small number of directions that contain most of the information in data. This is typically accomplished by performing a singular value decomposition (SVD) of \mathbf{D} and retaining the singular vectors corresponding to the largest singular values. A limitation of this procedure is that it is highly sensitive to outliers

in the data set. Recently there has been much work done to develop and analyze algorithms for PCA that are robust with respect to outliers. A common way to model outliers is as sparse vectors [1]. In seminal papers Candès et. al. and Chandrasekaran et. al. introduced the Principal Components Pursuit (PCP) convex program and proved its robustness to sparse outliers [2], [3]. Principal Components Pursuit poses the robust PCA problem as identifying a low rank matrix and a sparse matrix from their sum. The program is to minimize a weighted sum of the nuclear norm of the low rank matrix and the vector ℓ_1 norm of the sparse matrix subject to their sum being equal to the observed data matrix. The results in [4] improve upon those in [3]. Other methods such as [5] model the entire column vector as being either correct or an outlier. Some other works on the performance guarantees for batch robust PCA include [6], [7], and [8]. All of these methods require waiting until all of the data has been acquired before performing the optimization.

In this work we consider an online or recursive version of the robust PCA problem where we seek to separate vectors into low dimensional and sparse components as they arrive, using the previous estimates, rather than re-solving the entire problem at each time t . An application where this type of problem is useful is in video analysis [9]. Imagine a video sequence that has a distinct background and foreground. An example might be a surveillance camera where a person walks across the scene. If the background does not change very much, and the foreground is sparse (both practical assumptions), then separating the background and foreground can be viewed as a robust PCA problem. Sparse plus low rank decomposition can also be used to detect anomalies in network traffic patterns [10]. In all such an applications an online solution is desirable.

3.1.1 Paper Organization

This paper is organized in the following way. The signal model, assumptions, and main result are given in Section 3.2. Here we explain special cases of support change of \mathbf{x}_t that will satisfy our assumptions. Section 3.3 contains our most general support change model. Proofs of the support change results can be found in Appendix 3.B. A description of the Algorithm studied is in Section 3.4.

The main theorem (Theorem 3.2.15) is proved in Section 3.5. Proofs of the main lemmas used to prove the theorem are given next in Sections 3.6 and 3.7. Some preliminary and simple results used in the proofs are deferred to the Appendix. Appendix 3.A contains lemmas for exchanging the order of a double sum (Lemma 3.A.1), Cauchy-Schwarz for matrices (Lemma 3.A.3), and the matrix Azuma inequality from [11] and associated corollaries (Lemmas 3.A.8 - 3.A.10).

In Section 3.8 we introduce a more general subspace change model and a corresponding algorithm. A result analogous to Theorem 3.2.15 is proven in Section 3.9. Finally, in Section 3.10 we provide some simple simulation results that demonstrate the theoretical results.

3.1.2 Problem Definition

At time t we observe a vector $\mathbf{m}_t \in \mathbb{R}^n$ that is the sum of a vector from a slowly changing low-dimensional subspace ℓ_t and a sparse vector \mathbf{x}_t . So

$$\mathbf{m}_t = \ell_t + \mathbf{x}_t \quad \text{for } t = 1, 2, \dots, t_{\max},$$

with the possibility that $t_{\max} = \infty$. We model the low-dimensional ℓ_t 's as $\ell_t = \mathbf{P}_t \mathbf{a}_t$ for a basis matrix \mathbf{P}_t that is allowed to change slowly over time. Given an estimate of the initial subspace $\hat{\mathbf{P}}_{(0)}$, the goal is to obtain estimates $\hat{\mathbf{x}}_t$ and $\hat{\ell}_t$ at each time t and to periodically update the estimate of \mathbf{P}_t .

3.1.3 Contribution

The ReProCS algorithm for online sparse + low-rank matrix recovery (online robust PCA) was first introduced and analyzed in [12]. That paper contained a result that assumed certain properties of the algorithm estimates, and as such was only a partial correctness result. In Chapter 2, by building on proof techniques introduced in [12], a full correctness result was proven for the same ReProCS algorithm. Experimental evaluation of ReProCS is done in [14]. Here it is shown that with practical heuristics used to set its parameters (some of which are different from the algorithm parameters used in our correctness result), ReProCS has signifi-

cantly improved recovery performance compared to other recursive ([15, 16, 10]) and even batch methods ([2, 9, 16]) for many simulated and real video datasets.

A limitation of the result of Chapter 2 was that it assumed independence of the ℓ_t 's over time. For the video application, this means that the background images¹ are independent over time, and in most cases this is not a valid assumption. In this work, we are able to remove this assumption and replace it by a more realistic first order autoregressive (AR) model assumption on the ℓ_t 's. The ReProCS algorithm itself does not need knowledge of the AR model or its parameters. Under this model and only one extra assumption compared to Chapter 2, we can obtain a correctness result. The extra assumption needed is a bound on the ratio of the squared maximum value to the variance of any entry of \mathbf{a}_t (which will later be defined as the coefficients of ℓ_t with respect to an orthonormal basis). We show that as long as algorithm parameters are set appropriately, a good enough estimate of the initial subspace is available, slow subspace change holds, the subspaces are dense enough, and there is a certain amount of support change at least every so often, then the support can be exactly recovered with high probability, the sparse and low-rank matrix columns can be recovered with bounded and small error, and the subspace recovery error decays to a small value within a finite delay of a subspace change. Use of the AR model requires new proof techniques beyond what were used in Chapter 2. We need to use the matrix Azuma inequality from [11] instead of the matrix Hoeffding inequality from the same paper. Before applying the matrix Azuma inequality, we must also perform algebraic manipulation of sums so that the previous term on which we are conditioning is small. This is done in Lemma 3.6.3 which is proved using Lemma 3.A.6. Lemma 3.6.3 is used in the proof of Lemma 3.5.22 which is also significantly different.

A second contribution of this work is that we assume a subspace change model that allows for both addition of new directions and removal of existing directions from the subspace, and we also study a partly practical modification of the ReProCS with cluster-PCA (ReProCS-cPCA) algorithm introduced in [12]. ReProCS-cPCA improves upon ReProCS in that it also includes a subspace re-estimation step that allows removal of deleted directions from the subspace estimate. We obtain a correctness result for this algorithm under all the earlier assumptions

¹technically the background image minus a mean background image

and a clustering assumption on the eigenvalues of the covariance matrix of ℓ_t after the subspace change has stabilized. A key advantage of this result is that it significantly relaxes the denseness requirements and consequently needs a much looser upper bound on the rank-sparsity product compared with our result for ReProCS. The ReProCS result needs a bound that is tighter than what PCP needs (PCP is a batch method while ReProCS is online) but the bound needed by ReProCS-cPCA is comparable to that of PCP. In fact, ReProCS-cPCA does not need a bound on the rank of \mathbf{L} as long as the delay between subspace change times increases in proportion to $\log J$ where J is the total number of subspace change times in the entire sequence and J is known. The requirement that the length of time between subspace times increases with J is a consequence of the probabilistic signal model, and not the algorithm. Another way to interpret the result is that the probability of incorrect recovery increases only linearly with J . However this result has a significant limitation. Unlike the practical ReProCS algorithm, the ReProCS-cPCA algorithm is not fully automatic. It needs information about the eigenvalue clustering which at this point we cannot set automatically. (See the discussion in Section 3.8.3.)

To the best of our knowledge, this is among the first few works that provides a correctness result for an online (recursive) algorithm for sparse plus low-rank matrix recovery or equivalently for online robust PCA. As an easy corollary, we also have a result for online matrix completion. In this case, the support of \mathbf{x}_t is the set of missing entries, and this is known. Online algorithms are needed for real-time applications; even for offline applications, they are faster and need less storage compared to batch techniques. Moreover, online approaches can provide a natural way to exploit temporal dependencies in the dataset. In our case, we show that ReProCS uses slow subspace change to allow for significantly more correlated support sets of the sparse vectors than do the various results for PCP [2, 3, 17]. Partial results have been provided for online sparse plus low-rank matrix recovery in our earlier work [12], and also in later work by Feng et. al. [18]; however, all require an assumption on intermediate algorithm estimates. We discuss these and [19, 20] in Sec 3.2.6. There is some more recent work on online robust PCA algorithms and their experimental evaluation, e.g. [15], [21].

The new proof techniques developed in this and earlier works [12, 13] are needed because in our case, the error $\mathbf{e}_t = \ell_t - \hat{\ell}_t$ is correlated with the true data ℓ_t . This is an artifact of how

we obtain the estimates ℓ_t . Standard results for PCA such as those in [22] cannot be used, because they assume that the noise is independent of or uncorrelated with the noise-free data.

3.1.4 Notation

We use lowercase bold letters for vectors and capital bold letters for matrices. We use \mathbf{x}' for the transpose of a vector \mathbf{x} and similarly \mathbf{A}' for the transpose of a matrix \mathbf{A} . (Everything is real, so this is also the Hermitian adjoint). The 2-norm of a vector and the induced 2-norm of a matrix are denoted by $\|\cdot\|_2$. We refer to a matrix with orthonormal columns as a *basis matrix*. Notice that for a basis matrix \mathbf{P} , $\mathbf{P}'\mathbf{P} = \mathbf{I}$. For a set \mathcal{T} of integers, $|\mathcal{T}|$ denotes its cardinality. For a vector \mathbf{x} , $\mathbf{x}_{\mathcal{T}}$ is a smaller vector containing the entries of \mathbf{x} indexed by \mathcal{T} . Define $\mathbf{I}_{\mathcal{T}}$ to be an $n \times |\mathcal{T}|$ matrix of those columns of the identity matrix indexed by \mathcal{T} . Then let $\mathbf{A}_{\mathcal{T}} := \mathbf{A}\mathbf{I}_{\mathcal{T}}$. For matrices \mathbf{P} and \mathbf{Q} where the columns of \mathbf{Q} are a subset of the columns of \mathbf{P} , we will use the notation $\mathbf{P} \setminus \mathbf{Q}$ to mean the matrix of columns in \mathbf{P} and not in \mathbf{Q} . Using our column subscripting notation, for an $n \times r$ matrix \mathbf{P} , if $\mathbf{Q} = \mathbf{P}_{\mathcal{T}}$, then $\mathbf{P} \setminus \mathbf{Q} := \mathbf{P}_{[1, \dots, r] \setminus \mathcal{T}} = \mathbf{P}_{\overline{\mathcal{T}}}$.

For integers a and b , we write $a \bmod b$ for the remainder when a is divided by b . We use the interval notation $[a, b]$ to mean all of the integers between a and b , inclusive, and similarly for (a, b) etc. For a matrix \mathbf{A} , the restricted isometry constant (RIC) $\delta_s(\mathbf{A})$ is the smallest real number δ_s such that

$$(1 - \delta_s)\|\mathbf{x}\|_2^2 \leq \|\mathbf{A}\mathbf{x}\|_2^2 \leq (1 + \delta_s)\|\mathbf{x}\|_2^2$$

for all s -sparse vectors \mathbf{x} [23]. A vector \mathbf{x} is s -sparse if it has s or fewer non-zero entries. For Hermitian matrices \mathbf{A} and \mathbf{B} , the notation $\mathbf{A} \preceq \mathbf{B}$ means that $\mathbf{B} - \mathbf{A}$ is positive semi-definite. For an Hermitian matrix \mathbf{H} , $\mathbf{H} \stackrel{\text{EVD}}{=} \mathbf{U}\mathbf{\Lambda}\mathbf{U}'$ denotes its eigenvalue decomposition. Similarly for any matrix \mathbf{A} , $\mathbf{A} \stackrel{\text{SVD}}{=} \mathbf{U}\mathbf{\Sigma}\mathbf{V}'$ denotes its singular value decomposition.

For basis matrices \mathbf{P} and \mathbf{Q} , define

$$\text{dif}(\mathbf{P}, \mathbf{Q}) := \|(I - \mathbf{P}\mathbf{P}')\mathbf{Q}\|_2.$$

It is not difficult to show that this function is symmetric when \mathbf{P} and \mathbf{Q} are the same size. (See for example [12]). We will use the function dif as a measure of the error made when estimating

$\text{range}(\mathbf{Q})$ by $\text{range}(\mathbf{P})$. We use the hat to denote estimation of the column space of a matrix. So $\text{range}(\hat{\mathbf{P}})$ is an estimate of $\text{range}(\mathbf{P})$, but $\hat{\mathbf{P}}$ is not necessarily an entry-wise estimate of \mathbf{P} .

We will use \mathbb{P} for probability and \mathbb{E} for expectation. For a sequence of random variables $\mathbf{Z}_1, \dots, \mathbf{Z}_t$, the notation

$$\mathbb{E}_{t-1}[\mathbf{Z}_t] := \mathbb{E}[\mathbf{Z}_t | \mathbf{Z}_1, \dots, \mathbf{Z}_{t-1}]$$

and

$$\mathbb{E}_{t-1}[\mathbf{Z}_t | X] := \mathbb{E}[\mathbf{Z}_t | X, \mathbf{Z}_1, \dots, \mathbf{Z}_{t-1}]$$

For an event (set) Γ , we use $\bar{\Gamma}$ for the complement of Γ .

3.2 Model Assumptions, Main Result, and Discussion

3.2.1 Model on ℓ_t

Model 3.2.1.

1. Subspace Change Model for ℓ_t

Let t_j for $j = 1, \dots, J$ be the times at which the subspace changes. For the sake of notation, let $t_0 = 0$ and $t_{J+1} := t_{\max}$. We assume $\ell_t = \mathbf{P}_t \mathbf{a}_t$ for all $t = 1, \dots, t_{\max}$, and

$$\mathbf{P}_t = \begin{cases} [\mathbf{P}_{t-1} & \mathbf{P}_{t,\text{new}}] & \text{if } t = t_1 \text{ or } t_2 \text{ or } \dots t_J \\ \mathbf{P}_{t-1} & \text{otherwise} \end{cases} \quad (3.1)$$

where \mathbf{P}_t is a basis matrix for all t .

Let $r_j = \text{rank}(\mathbf{P}_{t_j})$ and $c_{j,\text{new}} = \text{rank}(\mathbf{P}_{t_j,\text{new}})$.

2. Assumptions and notation for \mathbf{a}_t . We assume the following model on \mathbf{a}_t :

$$\mathbf{a}_t = b\mathbf{a}_{t-1} + \boldsymbol{\nu}_t$$

for a scalar $b < 1$. Set $\mathbf{a}_0 = \mathbf{0}$. Also assume $\mathbb{E}[\boldsymbol{\nu}_t] = \mathbf{0}$, the $\boldsymbol{\nu}_t$ are mutually independent over t , bounded, and the matrix $\boldsymbol{\Lambda}_{\nu,t} := \text{Cov}(\boldsymbol{\nu}_t)$ is diagonal.

At $t = t_j$ the length of the vector \mathbf{a}_t increases from r_{j-1} to r_j . From (3.1) it is clear that for $t < t_j$, $\mathbf{P}_{t_j,\text{new}}' \ell_t = \mathbf{0}$. Therefore, for the autoregressive model, the last $(r_j - r_{j-1})$ entries of \mathbf{a}_t begin at 0.

Define

$$(a) \quad \gamma := \frac{\sup_t \|\boldsymbol{\nu}_t\|_\infty}{(1-b)} \quad (\text{Since } \boldsymbol{\nu}_t \text{ is bounded, } \gamma < \infty.)$$

$$(b) \quad \lambda^- := \frac{\inf_t \lambda_{\min}(\boldsymbol{\Lambda}_{\nu,t})}{1-b^2} \quad \text{and} \quad \lambda^+ := \frac{\sup_t \lambda_{\max}(\boldsymbol{\Lambda}_{\nu,t})}{1-b^2}$$

and assume that $0 < \lambda^- \leq \lambda^+ < \infty$.

By induction it is easy to see that

$$\|\mathbf{a}_t\|_\infty = \|b\mathbf{a}_{t-1} + \boldsymbol{\nu}_t\|_\infty \leq b\|\mathbf{a}_{t-1}\|_\infty + \|\boldsymbol{\nu}_t\|_\infty \leq b\gamma + (1-b)\gamma \leq \gamma.$$

Definition 3.2.2. Define

$$\mathbf{P}_{(j)} := \mathbf{P}_{t_j} \text{ for } j = 1, \dots, J$$

$$\mathbf{P}_{(j),*} := \mathbf{P}_{(j-1)}$$

$$\mathbf{P}_{(j),\text{new}} := \mathbf{P}_{t_j,\text{new}}$$

Because \mathbf{P}_t is a basis matrix, $\mathbf{P}_{(j),*} \perp \mathbf{P}_{(j),\text{new}}$. Define $r_J = \max_j \text{rank } \mathbf{P}_{(j)}$. and $c_{\text{new}} := \max_j \text{rank}(\mathbf{P}_{(j),\text{new}})$. Observe that

$$r_J = \text{rank}([\boldsymbol{\ell}_1, \dots, \boldsymbol{\ell}_{t_{\max}}]) \quad \text{and} \quad r_J \leq r_0 + Jc_{\text{new}} := r.$$

For $t \in [t_j, t_{j+1})$, $\boldsymbol{\ell}_t$ can be written as $\boldsymbol{\ell}_t = [\mathbf{P}_{(j),*} \quad \mathbf{P}_{(j),\text{new}}] \begin{bmatrix} \mathbf{a}_{t,*} \\ \mathbf{a}_{t,\text{new}} \end{bmatrix}$, where

$$\mathbf{a}_{t,*} := \mathbf{P}_{(j),*}' \boldsymbol{\ell}_t \quad \text{and} \quad \mathbf{a}_{t,\text{new}} := \mathbf{P}_{(j),\text{new}}' \boldsymbol{\ell}_t \quad (3.2)$$

Definition 3.2.3. As in (3.2), for $t \in [t_j, t_{j+1})$, define

$$\boldsymbol{\nu}_{t,*} := (\boldsymbol{\nu}_t)_{[1, r_j - c_{j,\text{new}}]} \quad \text{and} \quad \boldsymbol{\nu}_{t,\text{new}} := (\boldsymbol{\nu}_t)_{[r_j - c_{j,\text{new}} + 1, r_j]}$$

and define $\boldsymbol{\Lambda}_{\nu,t} := \text{Cov}(\boldsymbol{\nu}_t)$. Also define $\boldsymbol{\Lambda}_{a,t} := \text{Cov}(\mathbf{a}_t)$.

Observe that

$$\boldsymbol{\Lambda}_{a,t} = b^2 \boldsymbol{\Lambda}_{a,t-1} + \boldsymbol{\Lambda}_{\nu,t} \quad (3.3)$$

From the above equation, it is clear that $\boldsymbol{\Lambda}_{a,t}$ is also diagonal and

$$(1 - b^{2t})\lambda^- \leq \lambda_{\min}(\boldsymbol{\Lambda}_{a,t}) \leq \lambda_{\max}(\boldsymbol{\Lambda}_{a,t}) \leq (1 - b^{2t})\lambda^+.$$

Definition 3.2.4. Define

$$f := \frac{\lambda^+}{\lambda^-}$$

Notice that f is a bound on the condition number of $\mathbf{\Lambda}_{a,t}$ at any time t .

Definition 3.2.5. Define $\mathbf{\Lambda}_{a,t,*} := \text{Cov}(\mathbf{a}_{t,*})$ and $\mathbf{\Lambda}_{a,t,\text{new}} := \text{Cov}(\mathbf{a}_{t,\text{new}})$. Similarly define $\mathbf{\Lambda}_{\nu,t,*} := \text{Cov}(\boldsymbol{\nu}_{t,*})$ and $\mathbf{\Lambda}_{\nu,t,\text{new}} := \text{Cov}(\boldsymbol{\nu}_{t,\text{new}})$. Then

$$\mathbf{\Lambda}_{a,t} = \begin{bmatrix} \mathbf{\Lambda}_{a,t,*} & \mathbf{0} \\ \mathbf{0} & \mathbf{\Lambda}_{a,t,\text{new}} \end{bmatrix} \text{ and } \mathbf{\Lambda}_{\nu,t} = \begin{bmatrix} \mathbf{\Lambda}_{\nu,t,*} & \mathbf{0} \\ \mathbf{0} & \mathbf{\Lambda}_{\nu,t,\text{new}} \end{bmatrix}.$$

Definition 3.2.6. For an integer d , define

$$\begin{aligned} 1. \quad \gamma_{\text{new}} &:= \frac{\max_j \max_{t \in [t_j, t_j+d]} \|\boldsymbol{\nu}_{t,\text{new}}\|_{\infty}}{(1-b)} \\ 2. \quad \lambda_{\text{new}}^- &:= \min_j \min_{t \in [t_j, t_j+d]} \frac{\lambda_{\min}(\mathbf{\Lambda}_{\nu,t,\text{new}})}{1-b^2} \text{ and } \lambda_{\text{new}}^+ := \max_j \max_{t \in [t_j, t_j+d]} \frac{\lambda_{\max}(\mathbf{\Lambda}_{\nu,t,\text{new}})}{1-b^2} \end{aligned}$$

In the theorem we assume upper bounds on γ_{new} and λ_{new}^+ . This ensures that the projection of ℓ_t along the new directions is “small” for some time (d frames) after a subspace change.

Definition 3.2.7. Define

$$\eta := \frac{\gamma^2}{\lambda^+} \text{ and } \eta_{\text{new}} := \frac{\gamma_{\text{new}}^2}{\lambda_{\text{new}}^+}$$

Observe that if \mathbf{a}_t is a scalar random variable, then η is the ratio of the maximum absolute value (of \mathbf{a}_t) squared to the variance. For a continuous random variable uniformly distributed on the interval $[-\gamma, \gamma]$ one can easily compute $\eta = 3$. Moreover, if the i -th entry of \mathbf{a}_t is a continuous random variable that is uniformly distributed in $[-\gamma_i, \gamma_i]$ and if this is true for all times t , then $\gamma = \max_i \gamma_i$ and we still get $\eta = \eta_{\text{new}} = 3$.

Subspace Change Model with Deletions. The above simple model only assumes new directions get added to the subspace but nothing gets removed. We work with this model for notational simplicity. Our results in this section will also hold if it is replaced by the following more general model.

Model 3.2.8. Signal Model 3.2.1 holds with $\mathbf{P}_t = [(\mathbf{P}_{t-1}\mathbf{R}_t) \setminus \mathbf{P}_{t,\text{old}} \quad \mathbf{P}_{t,\text{new}}]$ at $t = t_j$'s. Here \mathbf{R}_t is an arbitrary rotation matrix. (By using a rotation matrix, we allow the removal of any direction in $\text{range}(\mathbf{P}_{t-1})$.)

In Section 3.8, we study a more general ReProCS algorithm called ReProCS-cPCA that includes a deletion (by subspace re-estimation) step. Our result for that algorithm is based on Signal Model 3.2.8. With one extra assumption, we are able to prove a stronger conclusion for ReProCS-cPCA.

3.2.2 Denseness coefficient

Below we give the definition of the denseness coefficient κ_s .

Definition 3.2.9. For a basis matrix \mathbf{P} , define $\kappa_s(\mathbf{P}) := \max_{|\mathcal{T}| \leq s} \|\mathbf{I}_{\mathcal{T}}' \mathbf{P}\|_2$.

As described in [12], κ_s is a measurement of the denseness of the vectors in the subspace $\text{range}(\mathbf{P})$. Notice that small κ_s means that the columns of \mathbf{P} are dense vectors. The reason for quantifying denseness using κ_s is the following lemma from [12].

Lemma 3.2.10. For a basis matrix \mathbf{P} , $\delta_s(\mathbf{I} - \mathbf{P}\mathbf{P}') = (\kappa_s(\mathbf{P}))^2$.

3.2.3 Model on \mathbf{x}_t

Let $\mathcal{T}_t := \{i : (\mathbf{x}_t)_i \neq 0\}$ be the support set of \mathbf{x}_t and let

$$s := \max_t |\mathcal{T}_t|$$

be the size of the largest support, and let

$$x_{\min} := \inf_t \min_{i \in \mathcal{T}_t} |(\mathbf{x}_t)_i|$$

denote the size of the smallest non-zero entry of any \mathbf{x}_t .

In order to prove our result, we require that the supports of \mathbf{x}_t be sufficiently different. Below we give two possible models that will imply the most general conditions required (given in Section 3.3).

Model 3.2.11. Let $\varrho \geq 1$ be a scalar. Suppose that the support of \mathbf{x}_t is of fixed size less than or equal to s , consists of consecutive indices, and moves down the vector at least every β time instants. Moreover, each time it moves, it moves by at least $\frac{s}{\varrho}$ indices and at most $\varrho_2 s$ indices. Stated differently, the support of \mathbf{x}_t remains the same for no more than β time instants, and when it moves, it moves by no fewer than $\frac{s}{\varrho}$ indices and no more than $\varrho_2 s$ indices.

Model 3.2.12. Suppose that the support of \mathbf{x}_t consists of s or fewer consecutive indices and moves down the vector by between 1 and m indices at every time t .

Figure 3.5 on page 98 illustrates the above support change models.

Remark 3.2.13. For both of the above models, when the support reaches the bottom of the vector, we assume that it starts over at 1. This models a moving 1D object of length s or less that enters the scene and eventually walks out, and then another object of length s or less may come in. The requirement of consecutive indices and downward (as opposed to upward) motion are done for simplicity and ease of understanding. Our results still hold under permutations (relabeling) of the indices. We could also make a small modification and assume that the object is reflected back up (down) when it reaches the bottom (top). See Remark 3.3.6.

Remark 3.2.14. Nothing in our most general support change model or our algorithm requires only one object in the support of \mathbf{x}_t . The models above are simple examples that capture the intuition of our general model.

3.2.4 Main Result

Theorem 3.2.15. Consider Algorithm 3. Pick a ζ that satisfies

$$\zeta \leq \min \left\{ \frac{10^{-4}}{r^2}, \frac{1.5 \times 10^{-4}}{r^2 f}, \frac{1}{r^3 \gamma^2}, \frac{0.01 \lambda^-}{b^2 r^3 \gamma^2} \right\}.$$

Suppose

1. $\|(\mathbf{I} - \hat{\mathbf{P}}_{(0)} \hat{\mathbf{P}}_{(0)}') \mathbf{P}_{(0)}\|_2 \leq r_0 \zeta$;
2. The algorithm parameters are set as:
 - $\text{thresh} = \frac{\lambda^-}{2}$;
 - $K = \left\lceil \frac{\log(0.16 c_{\text{new}} \zeta)}{\log(0.4)} \right\rceil$;
 - $\xi = \sqrt{c_{\text{new}}} \gamma_{\text{new}} + (\sqrt{r} + \sqrt{c_{\text{new}}}) \sqrt{\zeta}$;
 - $\omega = 7\xi$;

- $\alpha = C(\log(35KJ) + 11 \log(n))$ for a constant $C \geq C_{\text{add}}$ with

$$C_{\text{add}} := 20^2 \cdot 8 \cdot 48^2 \frac{(\phi^+ \xi)^4}{(c_{\text{new}} \zeta \lambda^-)^2}; \quad (3.4)$$

3. Signal Model 3.2.1 on ℓ_t holds with $b \leq 0.1$ and

(a) $t_{j+1} - t_j > d \geq (K+2)\alpha$ for all j i.e. the delay between change times is “large”;

(b) $\mathbf{a}_{t,\text{new}}$ is “small”

$$i. \sqrt{c_{\text{new}}} \gamma_{\text{new}} + (\sqrt{r} + \sqrt{c_{\text{new}}}) \sqrt{\zeta} \leq \frac{x_{\min}}{14};$$

$$ii. \lambda_{\text{new}}^+ \leq \sqrt{2} \lambda^-;$$

$$iii. b^2 c_{\text{new}} \eta_{\text{new}} \lambda_{\text{new}}^+ \leq 0.5 \lambda^- \text{ (because } b \leq 0.1, \text{ this will be satisfied if } c_{\text{new}} \eta_{\text{new}} \lambda_{\text{new}}^+ \leq 50 \lambda^-);$$

4. The support of \mathbf{x}_t changes enough so that for the α chosen above, Signal Model 3.2.11 holds with $\beta = h^+ \alpha$ and

$$\lceil \varrho \rceil^2 h^+ \leq 0.0025, \text{ and } \varrho_2 s \alpha \leq n$$

or Signal Model 3.2.12 holds with $s \leq (6 \times 10^{-4})\alpha$ and $\alpha \leq \frac{n}{m}$.

5. The low dimensional subspace is dense such that

$$\kappa_{2s}(\mathbf{P}_{(J)}) \leq 0.3 \text{ and } \max_j \kappa_{2s}(\mathbf{P}_{(j),\text{new}}) \leq 0.02.$$

Then, with probability at least $1 - n^{-10}$, at all times t

1. The support of \mathbf{x}_t is recovered exactly, i.e. $\hat{\mathcal{T}}_t = \mathcal{T}_t$;
2. The estimate of the subspace change time satisfies $t_j \leq \hat{t}_j \leq t_j + 2\alpha$, for $j = 1, \dots, J$;
3. The estimate of the number of new directions is correct, i.e. $\hat{c}_{j,\text{new},k} = c_{j,\text{new}}$ for $j = 1, \dots, J$ and $k = 1, \dots, K$;

4. The recovery error satisfies:

$$\|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2 \leq \begin{cases} 1.2 (\sqrt{\zeta} + \sqrt{c_{\text{new}}\gamma_{\text{new}}}) & t \in [t_j, \hat{t}_j] \\ 1.2 (1.84\sqrt{\zeta} + (0.4)^{k-1}\sqrt{c_{\text{new}}\gamma_{\text{new}}}) & t \in [\hat{t}_j + (k-1)\alpha, \hat{t}_j + k\alpha - 1], \\ & k = 1, 2, \dots, K \\ 2.4\sqrt{\zeta} & t \in [\hat{t}_j + K\alpha, t_{j+1} - 1]; \end{cases}$$

5. The subspace error $\text{SE}_t := \|(\mathbf{I} - \hat{\mathbf{P}}_t \hat{\mathbf{P}}_t') \mathbf{P}_t\|_2$ satisfies:

$$\text{SE}_t \leq \begin{cases} 1 & t \in [t_j, \hat{t}_j] \\ 10^{-2}\sqrt{\zeta} + 0.4^{k-1} & t \in [\hat{t}_j + (k-1)\alpha, \hat{t}_j + k\alpha - 1], \quad k = 1, 2, \dots, K \\ 10^{-2}\sqrt{\zeta} & t \in [\hat{t}_j + K\alpha, t_{j+1} - 1]. \end{cases}$$

Remark 3.2.16. When $b = 0$ (\mathbf{a}_t 's are independent), the bound $\zeta \leq \frac{0.01\lambda^-}{b^2 r^3 \gamma^2}$ gets removed because $1/b^2 \rightarrow \infty$; and the third requirement in condition 3b) gets removed. Thus, when $b = 0$, the above theorem is the same as the main result in Chapter 2 even though it provides guarantees for a practical version of the ReProCS algorithm studied in Chapter 2.

3.2.5 Random Support Change

We give here a commonly used Bernoulli-Gaussian motion model that satisfies our support change assumptions with high probability.

Model 3.2.17 (Taken from Chapter 2). Consider one-dimensional motion of the support of \mathbf{x}_t , and let \tilde{o}_t be its center at time t . Suppose that the support moves according to the model

$$o_t = o_{t-1} + \theta_t \left(1.1 \frac{s}{\varrho} + \mu_t \right) \quad \text{and} \\ \tilde{o}_t = o_t \mod n$$

where μ_t is Gaussian $\mathcal{N}(0, \sigma^2)$ and θ_t is a Bernoulli random variable that takes the value 1 with probability q and 0 with probability $1 - q$, and $\varrho \geq 1$ is a constant. Taking the modulus with respect to n describes the process of the support starting over at 1 when it reaches n . Assume that $\{\mu_t\}$, $\{\theta_t\}$ are mutually independent and independent of $\{\nu_t\}$ for $t = 1, \dots, t_{\max}$.

Lemma 3.2.18 (Taken from Chapter 2). *Signal Model 3.2.17 satisfies the assumptions of Signal Model 3.2.11 with $\varrho_2 = \frac{1.2}{\varrho}$ with probability at least $1 - n^{-10}$ if $\sigma^2 \leq \frac{s^2}{4000\varrho^2 \log(n)}$, $q \geq 1 - \left(\frac{n^{-10}}{2(t_{\max} + \alpha)} \right)^{\frac{1}{\beta}}$, and $t_{\max} \leq n^{10}$.*

Proof. To prove the above lemma, we will be done if we can show

1. the support changes at least once every β instants with probability greater than $1 - \frac{n^{-10}}{2}$;
2. when it changes, the support moves by at least $\frac{s}{\varrho}$ and not more than $1.2\frac{s}{\varrho}$ indices with probability greater than $1 - \frac{n^{-10}}{2}$.

Item 1) follows using simple arguments for Bernoulli random variables [24] while item 2) follows using a Gaussian tail bound. The complete proof is in Appendix 3.B. \square

Corollary 3.2.19. *Consider Theorem 3.2.15 with condition 4) replaced by Signal Model 3.2.17 with*

$$(i) \quad \frac{1.2s\alpha}{\varrho} \leq n;$$

$$(ii) \quad t_{\max} \leq n^{10}$$

$$(iii) \quad \sigma^2 \leq \frac{s^2}{4000\varrho^2 \log(n)};$$

$$(iv) \quad q \geq 1 - \left(\frac{n^{-10}}{2(t_{\max} + \alpha)} \right)^{\frac{1}{\beta}} ; \text{ for a } \beta = \frac{(2.4 \times 10^{-3})\alpha}{\lceil \varrho \rceil^2}.$$

Then all its conclusions will hold with probability greater than or equal to $1 - 2n^{-10}$.

Corollary 3.2.19 follows by combining Lemma 3.2.18 with Theorem 3.2.15.

3.2.6 Discussion

The above result needs an accurate estimate of the initial subspace, a slow subspace change assumption, a support change assumption, and a denseness assumption. If a short sequence of background only training data is available (which is often true of surveillance video), then the initial subspace estimate is easy to obtain by ordinary PCA. Otherwise, one could use a

batch method for robust PCA on an initial sequence to set $\hat{\mathbf{P}}_{(0)}$ before starting the ReProCS algorithm.

Consider the subspace change model. This model (along with the bound on γ_{new} from the theorem) assumes that after a subspace change, $\|\mathbf{a}_{t,\text{new}}\|_\infty$ and therefore also $\|\mathbf{\Lambda}_{a,t,\text{new}}\|_2$ are initially small. After $t_j + d$, $\|\mathbf{a}_{t,\text{new}}\|_\infty$ can be as large as $\gamma = \max_t \|\mathbf{a}_t\|_\infty$, and the eigenvalues of $\mathbf{\Lambda}_{a,t,\text{new}}$ can increase up to λ^+ either immediately or gradually. Thus a new direction added at time t_j can have magnitude as large as γ and variance as large as λ^+ by $t_j + d$. Since we assume that $t_j + d \leq t_{j+1}$, this will occur before the next subspace change time. As demonstrated in [12], such a slow subspace change assumption is valid for backgrounds in real video sequences.

Consider the support change models. Both Models 3.2.11 and 3.2.12 are valid and commonly used models for foreground object motion in videos. Of course these are only special cases. Nothing prevents multiple moving objects (see Section 3.3.1). In the rest of this discussion we use Model 3.2.11. If we assume Model 3.2.11, our result requires $s \leq \frac{n}{2\alpha}$. If $J \leq C_1 n$ for some constant C_1 , then using the definition of α , this bound holds if $s \leq C_2 \frac{n}{\log n}$, for a constant C_2 . Thus, this model allows $s \in \mathcal{O}(\frac{n}{\log n})$ and $r = r_0 + Jc_{\text{new}} \in \mathcal{O}(n)$. As we explain next, our denseness assumption restricts the requirement on r to $r \in \mathcal{O}(\log n)$.

Consider denseness. The way κ_s is defined, our denseness assumption simultaneously places restrictions on denseness of ℓ_t , and on r and s . As done in [2], we could assume $\kappa_1(\mathbf{P}_{(J)}) \leq \sqrt{\frac{\mu r}{n}}$ and $\kappa_1(\mathbf{P}_{(j),\text{new}}) \leq \sqrt{\frac{\tilde{\mu} c_{\text{new}}}{n}}$ where μ and $\tilde{\mu}$ take any value between 1 and $\frac{n}{r}$. It is easy to show that $\kappa_s(\mathbf{P}) \leq \sqrt{s} \kappa_1(\mathbf{P})$ [12]. Thus if

$$\frac{2sr}{n} \leq \mu^{-1}(0.3)^2, \text{ and } \frac{2sc_{\text{new}}}{n} \leq \tilde{\mu}^{-1}(0.02)^2,$$

then our assumption of $\kappa_{2s}(\mathbf{P}_{(J)}) \leq 0.3$ and $\kappa_{2s}(\mathbf{P}_{(j),\text{new}}) \leq 0.02$ will be satisfied. Since we require $s \in \mathcal{O}(\frac{n}{\log n})$, this means we can allow $r \in \mathcal{O}(\log n)$ to satisfy the above.

Comparison with other work.

Let $\mathbf{L} = [\ell_1, \ell_2, \dots, \ell_{t_{\text{max}}}]$ and $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{t_{\text{max}}}]$. Define $r_{\text{mat}} := \text{rank}(\mathbf{L})$ and $s_{\text{mat}} := |\text{support}(\mathbf{X})|$. For our model, $s_{\text{mat}} = st_{\text{max}}$, and $r_{\text{mat}} = r_J \leq r$. From the above discussion, we see that the ReProCS result allows

$$s_{\text{mat}} \in \mathcal{O}\left(\frac{nt_{\text{max}}}{\log n}\right) \text{ and } r_{\text{mat}} \in \mathcal{O}(\log n). \quad (3.5)$$

The above requirement on s_{mat} and r_{mat} is stronger than that used by [2] (which studies the batch approach PCP). The result in [2] allows

$$s_{\text{mat}} \in \mathcal{O}(nt_{\text{max}}) \text{ and } r_{\text{mat}} \in \mathcal{O}\left(\frac{n}{(\log n)^2}\right).$$

But, up to differences in the constants, (3.5) is the same as the requirement found in [25] (which also studies the PCP program and is an improvement over [3]), except that [25] does not need specific bounds on s_{mat} and r_{mat} ; it only requires $r_{\text{mat}}s_{\text{mat}} \in \mathcal{O}(nt_{\text{max}})$. The comparison is not direct though because our result does not need denseness of the right singular vectors of \mathbf{L} or a bound on the vector infinity norm of \mathbf{UV}' , while [2, 3], and [25] do. Here $\mathbf{L} \stackrel{\text{SVD}}{=} \mathbf{U}\mathbf{\Sigma}\mathbf{V}'$. The reason for our stronger requirement on $s_{\text{mat}}r_{\text{mat}}$ is because we study an online algorithm, ReProCS, that recovers the sparse vector \mathbf{x}_t at each time t rather than in a batch or a piecewise batch fashion. Because of this, the sparse recovery step does not use the low dimensionality of the new (and still unestimated) subspace.

Because we only require that the support changes after a given maximum allowed duration, it can be constant for a certain period of time (Model 3.2.11), or it can change only a little at each time (Model 3.2.12). This is a substantially weaker assumption than the independent or uniformly random supports required by [2] and [20]. As we explain in Chapter 2, if we consider the whole matrix \mathbf{X} , then at most $\frac{t_{\text{max}}}{5000}$ non-zero entries per row are allowed by our result. Thus, for $r_{\text{mat}} > 5000$, this also is a significant improvement over [25] which requires at most $\frac{t_{\text{max}}}{r}$ non-zero entries per row. Therefore, an important advantage of our result is that it allows for highly correlated support sets of \mathbf{x}_t , which is important for applications such as video surveillance that involve one or more moving foreground objects or persons forming the sparse vector \mathbf{x}_t .

Now consider works that also use initial subspace knowledge. Our result improves upon [12]’s results by removing the denseness requirements on $(\mathbf{I} - \mathbf{P}_{(j),\text{new}}\mathbf{P}_{(j),\text{new}}')\hat{\mathbf{P}}_t$ and $(\mathbf{I} - \hat{\mathbf{P}}_{(j),*}\hat{\mathbf{P}}_{(j),*}' - \hat{\mathbf{P}}_{t,\text{new}}\hat{\mathbf{P}}_{t,\text{new}}')\mathbf{P}_{(j),\text{new}}$ and thus provides a complete correctness result. It also improves on the results of Chapter 2 as explained earlier by studying a fully automatic version of ReProCS and assuming an autoregressive model on the ℓ_t ’s. In [18], Feng et. al. propose a method for online robust PCA and prove a partial result for their algorithm. The approach is

to reformulate the PCP program and use this reformulation to develop a recursive algorithm that converges asymptotically to the solution of PCP as long as the basis estimate $\hat{\mathbf{P}}_t$ is full rank at each time t . Since this result assumes something about the algorithm estimates, it is only a *partial* result. Another work of Feng et. al. [19] on online robust PCA does not model the outlier as a sparse vector but defines anything that is far from the data subspace as an outlier. Another recent work that uses knowledge of the initial subspace estimate is modified-PCP [20]. However, like PCP, this also needs uniformly random supports. Moreover it is a piecewise batch approach.

Limitations. One limitation of this work is that we do not prove exact recovery of \mathbf{x}_t or ℓ_t .

In order to set the algorithm parameters as assumed in Theorem 3.2.15 one would need knowledge of the model parameters γ , γ_{new} , r_0 , c_{new} , λ^+ , λ^- , and J . This is a rather impractical assumption; however, as shown below in Corollary 3.2.20, if an initial sequence of just ℓ_t 's is available, then this limitation can be handled with some additional assumptions.

Suppose that an initial sequence of just ℓ_t 's without any sparse corruptions is available. In the video surveillance application, this would correspond to having a short sequence of background only frames. Such a sequence can be used to obtain $\hat{\mathbf{P}}_{(0)}$ (as the eigenvectors corresponding to non-zero eigenvalues of $\sum_{t=1}^{t_{\text{train}}} \ell_t \ell_t'$) and to estimate several of the required model parameters. Let t_{train} be the length of the sequence and assume that \mathbf{P}_t is constant for the duration. Define

$$\begin{aligned}\lambda^+ &:= \lambda_{\max} \left(\frac{1}{t_{\text{train}}} \sum_{t=1}^{t_{\text{train}}} \ell_t \ell_t' \right) \\ \lambda^- &:= \min_{\lambda_i \neq 0} \lambda_i \left(\frac{1}{t_{\text{train}}} \sum_{t=1}^{t_{\text{train}}} \ell_t \ell_t' \right) \\ r_0 &:= \text{rank} \left(\sum_{t=1}^{t_{\text{train}}} \ell_t \ell_t' \right) \\ \gamma &:= \max_{t \in [1, t_{\text{train}}]} \|\mathbf{P}_t' \ell_t\|_{\infty}.\end{aligned}$$

Assume that

1. $c_{\text{new}} = 0.1r_0$ and $\gamma_{\text{new}} = 0.1\gamma$

2. $\frac{\lambda_{\max}(\mathbf{\Lambda}_{\nu,t})}{1-b^2} \leq \lambda^+$ for all t
3. $\frac{\lambda_{\min}(\mathbf{\Lambda}_{\nu,t})}{1-b^2} \geq \lambda^-$ for all t
4. $c_{j,\text{new}} \leq c_{\text{new}} = 0.1r_0$ for all j
5. $\frac{\|\nu_t\|_\infty}{1-b} \leq \gamma$ for all t
6. $\frac{\|\nu_{t,\text{new}}\|_\infty}{1-b} \leq \gamma_{\text{new}} = 0.1\gamma$ for $t \in [t_j, t_j + d]$ for all j .

Under these assumptions, the only remaining unknown parameter is J . In Theorem 3.2.15, knowledge of J is only used to control the probability with which the result holds. Therefore, we can state the following corollary.

Corollary 3.2.20. *Suppose that the above assumptions hold. In Theorem 3.2.15, replace the J in the expression for α with a 1, and suppose that all of the remaining assumptions of Theorem 3.2.15 are satisfied. Then all of the conclusions will also hold with probability at least $1 - Jn^{-10}$.*

Another limitation is that Signal Model 3.2.1 only allows for additions to the subspace. The more realistic model (Signal Model 3.2.8) also allows removals from the prior subspace. Nothing in our proof changes if we incorporate removals into the signal model, and we have the following corollary.

Corollary 3.2.21. *Theorem 3.2.15 also holds with Signal Model 3.2.1 on ℓ_t replaced by Signal Model 3.2.8.*

Intuitively, Signal Model 3.2.8 is a special case of Signal Model 3.2.1 because Model 3.2.8 restricts the subspace where the ℓ_t can lie. However, because of the minimum variance requirement (λ^-) of Model 3.2.1, we need technically need a separate proof. The proof remains the same; the only difference is that instead of estimating $\text{span}(\mathbf{P}_{(j)})$, the algorithm will maintain an estimate of $\text{span}([\mathbf{P}_{(0)}, \mathbf{P}_{(1),\text{new}}, \dots, \mathbf{P}_{(j),\text{new}}])$. Notice that when there are no directions deleted, these are equivalent. Another way of describing this is that although directions are deleted from the subspace in Signal Model 3.2.8, Algorithm 3 does not delete anything from its subspace estimate. In Section 3.8 we introduce another algorithm (Algorithm 4) that does

remove old directions from its subspace estimate. A result similar to Theorem 3.2.15 is proven for this algorithm as well.

A fundamental limitation of our analysis approach is the assumption that there is a significant delay between times when the low-dimensional subspace changes. A more realistic model would allow the subspace to change more frequently.

Another limitation of the ReProCS algorithm is that it does not use the fact that the ‘noise’ seen by the sparse recovery step has a low dimensional structure. The modified PCP program from [20] uses this structure, but because of the piecewise batch approach cannot handle highly correlated supports of the sparse component like ReProCS (see the simulations in Section 3.10). The sparse recovery also requires a lower bound on x_{\min} to recover the support.

3.3 Most General Support Change Model and a Key Lemma

We give here the most general support change model under which our result holds. We will show that this includes the preceding signal models as special cases.

3.3.1 Most General Support Change Model

The model given below is a simple generalization of the support change model used in Chapter 2.

Model 3.3.1. *Let ρ be a positive integer. Split $[1, t_{\max}]$ into intervals of length α . Use $\mathcal{J}_u := [(u-1)\alpha + 1, u\alpha]$ to denote the u -th interval. For a given interval, \mathcal{J}_u , let $\mathcal{T}_{(i),u}$ for $i = 1, \dots, l_u$ be mutually disjoint subsets of $\{1, \dots, n\}$ such that for every $t \in \mathcal{J}_u$,*

$$\mathcal{T}_t \subseteq \mathcal{T}_{(i),u} \cup \mathcal{T}_{(i+1),u} \cup \dots \cup \mathcal{T}_{(i+\rho-1),u} \quad \text{for some } i. \quad (3.6)$$

For these $\mathcal{T}_{(i),u}$ ’s define

$$h \left(\alpha; \left\{ \mathcal{T}_{(i),u} \right\}_{\substack{u=1, \dots, \lceil \frac{t_{\max}}{\alpha} \rceil \\ i=1, \dots, l_u}} \right) := \max_{u=1, \dots, \lceil \frac{t_{\max}}{\alpha} \rceil} \max_i \left| \left\{ t \in \mathcal{J}_u : \mathcal{T}_t \subseteq \mathcal{T}_{(i),u} \cup \mathcal{T}_{(i+1),u} \cup \dots \cup \mathcal{T}_{(i+\rho-1),u} \right\} \right| \quad (3.7)$$

Now define $h^*(\alpha)$ which takes the minimum over all choices of $\mathcal{T}_{(i),u}$

$$h^*(\alpha) := \min_{\substack{\text{All choices of} \\ \text{mutually disjoint} \\ \mathcal{T}_{(i),u} \text{ satisfying (3.6)}}} h \left(\alpha; \{ \mathcal{T}_{(i),u} \}_{u=1, \dots, \lceil \frac{t_{\max}}{l_u^\alpha} \rceil} \right) \quad (3.8)$$

Assume that $h^*(\alpha) \leq h^+ \alpha$ for an $h^+ < 1$.

In the above model, $h^*(\alpha)$ provides a bound on how long the support of \mathbf{x}_t remains in a given area during an interval \mathcal{J}_u .

Notice that (3.6) can always be trivially satisfied by choosing $l_u = 1$ and $\mathcal{T}_{(1),u} = \{1, \dots, n\}$, but this will give $h(\alpha; \{\mathcal{T}_{(i),u}\}) = \alpha$ and hence is not a good choice. This is why we take a minimum over all choices.

The following corollary is the most general form of our result.

Theorem 3.3.2. *Suppose that the conditions of Theorem 3.2.15 hold, but instead of 4), assume that Signal Model 3.3.1 holds with $\rho^2 h^+ \leq .0024$. Then all its conclusions will hold.*

In Sec. 3.3.3, we show that Signal Models 3.2.11 and 3.2.12 are special cases of Signal Model 3.3.1 and hence Theorem 3.2.15 is actually a corollary of Theorem 3.3.2. The rest of the paper will prove Theorem 3.3.2.

3.3.2 Key Lemma

Under Signal Model 3.3.1, we can obtain the following lemma for a matrix $\mathbf{M} = \sum_t \mathbf{I}_{\mathcal{T}_t} \mathbf{A}_t \mathbf{I}_{\mathcal{T}_t}'$ (i.e. a sum of matrices supported on rows and columns indexed by \mathcal{T}_t). This lemma tells us that, because of the support change model, such a matrix is actually block banded: for $\varrho = 1$ it is block diagonal, for $\varrho = 2$ it is block tri-diagonal, and so on. Hence its 2-norm is much smaller than the sum of the norms of the individual matrices. As we will see later, the error, \mathbf{e}_t , in recovering \mathbf{x}_t (which is equal to the error in recovering $\boldsymbol{\ell}_t$) at times t , is supported on \mathcal{T}_t . As a result the matrix $\sum_t \mathbf{e}_t \mathbf{e}_t'$ has this form. Moreover, certain matrices obtained when bounding $\sum_t \boldsymbol{\ell}_t \mathbf{e}_t'$ also have this form.

Lemma 3.3.3 (Taken from Chapter 2). *Let $s_t = |\mathcal{T}_t|$. Consider a sequence of $s_t \times s_t$ symmetric positive-semidefinite matrices \mathbf{A}_t such that $\|\mathbf{A}_t\|_2 \leq \sigma^+$ for all t . Assume that the \mathcal{T}_t obey Signal*

Model 3.3.1 and the assumptions of Theorem 3.3.2. Let $\mathbf{M} = \sum_{t \in \mathcal{J}_u} \mathbf{I}_{\mathcal{T}_t} \mathbf{A}_t \mathbf{I}_{\mathcal{T}_t}'$ be an $n \times n$ matrix (\mathbf{I} is an $n \times n$ identity matrix). Then

$$\|\mathbf{M}\|_2 \leq \rho^2 h^+ \alpha \sigma^+ \leq 0.0024 \sigma^+ \alpha$$

Proof. This lemma is proved in Chapter 2. □

3.3.3 Unifying the signal models

Lemma 3.3.4. *Suppose that Signal Model 3.2.11 and the conditions assumed in Theorem 3.2.15 hold. Recall that this means that the support of \mathbf{x}_t is of a constant size s and moves by at least $\frac{s}{\varrho}$ indices, and at most $\varrho_2 s$ indices, at least every β time instants, $\beta = h^+ \alpha$, $\lceil \varrho \rceil^2 h^+ \leq 0.0024$ and $\varrho_2 s \alpha \leq n$. Then Signal Model 3.3.1 holds with $\rho = \lceil \varrho \rceil$ and $h^+ = \beta / \alpha$.*

The proof uses the same technique as in Chapter 2 and is given in Appendix 3.B, also see Figure 3.5.

Lemma 3.3.5. *Suppose that Signal Model 3.2.12 and the conditions assumed in Theorem 3.2.15 hold. Recall that this means that the support of \mathbf{x}_t moves by between 1 and m indices at every time t , $s \leq (3 \times 10^{-4}) \alpha$ and $\alpha \leq \frac{n}{m}$. Then Signal Model 3.3.1 is satisfied with $\rho = 2$ and $h^+ = s / \alpha$.*

The proof is simple and is given in Appendix 3.B, also see Figure 3.5.

Remark 3.3.6. *In Signal Models 3.2.11 and 3.2.12, if we replace the assumption that the support restarts at 1 when it reaches n , and instead assume that it is reflected back up the vector, then a form of Lemma 3.3.3 still holds, but the conclusion becomes $\|\mathbf{M}\|_2 \leq 2\rho^2 h^+ \sigma^+ \alpha$. The statement of Theorems 3.2.15 and 3.3.2 would then need the tighter bound on $\rho^2 h^+ \leq .0012$.*

In Lemma 3.2.18, we have already shown that Signal Model 3.2.17 is a special case of Signal Model 3.2.11. Above we have finished showing that both Signal Model 3.2.11 and 3.2.12 are special cases of Signal Model 3.3.1. In the rest of the paper, we assume Signal Model 3.3.1 and use Lemma 3.3.3 to prove Theorem 3.3.2. Theorem 3.2.15 follows as a special case of Theorem 3.3.2, which is the most general form of our result.

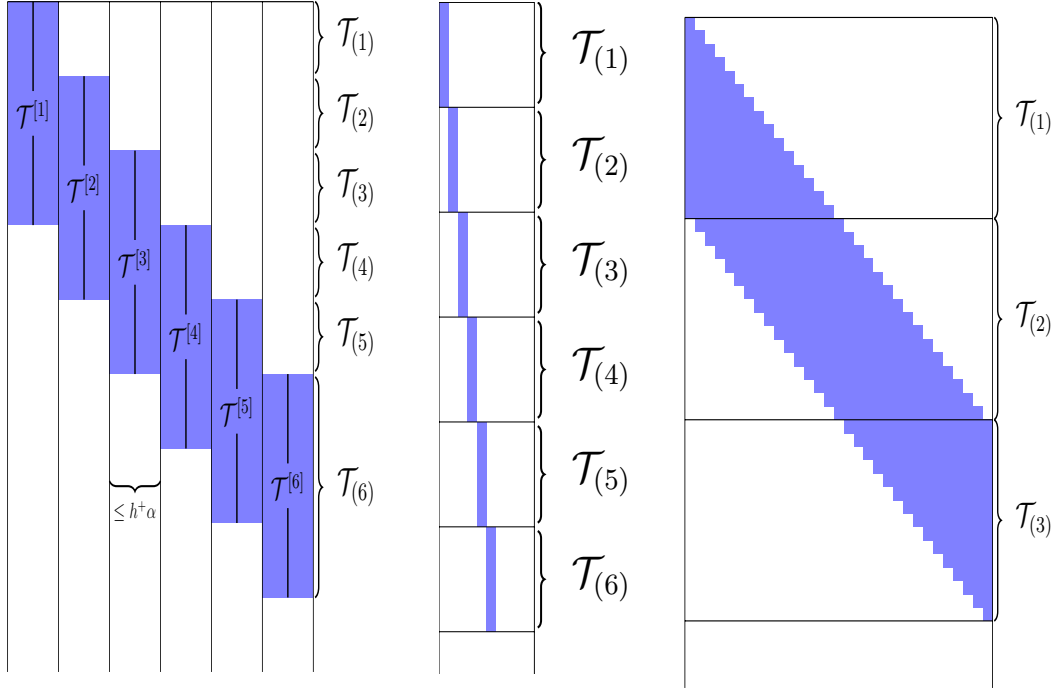
Figure 3.2 Signal Model 3.3.1 with $\varrho = 3$ Figure 3.3 Disjoint Supports (Signal Model 3.3.1 with $\varrho = 1$)

Figure 3.4 Signal Model 3.2.12

Figure 3.5: Support Change Models

3.4 The Automatic ReProCS Algorithm

The ReProCS algorithm was introduced in [12]. A more practical version including heuristics for setting the parameters was given in [14]. The basic idea of ReProCS is as follows. Given an accurate estimate of the subspace where the ℓ_t 's lie, projecting the measurement $\mathbf{m}_t = \mathbf{x}_t + \ell_t$ onto the orthogonal complement of the estimated subspace will nullify most of ℓ_t . The denseness of ℓ_t implies that this projection will have small RIC (Lemma 3.2.10) so the sparse recovery step will produce an accurate estimate $\hat{\mathbf{x}}_t$. Then, subtraction also gives a good estimate $\hat{\ell}_t = \mathbf{m}_t - \hat{\mathbf{x}}_t$. Using these $\hat{\ell}_t$, the algorithm successively updates the subspace estimate by a modification of the standard PCA procedure, which we call projection PCA.

Algorithm 3 is the same algorithm given in Chapter 2. The detailed description is given in Section 2.4

Algorithm 3 Recursive Projected CS (ReProCS) [12]

Parameters: $\xi, \omega, \alpha, K, \text{thresh}$

Input: \mathbf{m}_t , *Output:* $\hat{\mathbf{x}}_t, \hat{\boldsymbol{\ell}}_t, \hat{\mathbf{P}}_t, \hat{t}_j, \hat{c}_{j,\text{new},k}$

Notes: using $\hat{\mathbf{P}}_{t,*}, \hat{\mathbf{P}}_{t,\text{new}}, \hat{\mathbf{P}}_t = [\hat{\mathbf{P}}_{t,*} \hat{\mathbf{P}}_{t,\text{new}}]$

Set $\hat{\mathbf{P}}_{t,*} \leftarrow \hat{\mathbf{P}}_{(0)}, \hat{\mathbf{P}}_{t,\text{new}} \leftarrow [\cdot], j \leftarrow 0, \text{phase} \leftarrow \text{detect}$

For every $t > 0$, do the following:

1. Estimate \mathcal{T}_t and \mathbf{x}_t via Projected Compressed Sensing:
 - (a) Projection: set $\boldsymbol{\Phi}_t \leftarrow \mathbf{I} - \hat{\mathbf{P}}_{t-1} \hat{\mathbf{P}}_{t-1}'$, compute $\mathbf{y}_t \leftarrow \boldsymbol{\Phi}_t \mathbf{m}_t$
 - (b) Sparse Recovery: compute $\hat{\mathbf{x}}_{t,\text{cs}}$ as the solution of $\min_{\mathbf{x}} \|\mathbf{x}\|_1$ s.t. $\|\mathbf{y}_t - \boldsymbol{\Phi}_t \mathbf{x}\|_2 \leq \xi$
 - (c) Support Estimate: compute $\hat{\mathcal{T}}_t = \{i : |(\hat{\mathbf{x}}_{t,\text{cs}})_i| > \omega\}$
 - (d) LS Estimate of \mathbf{x}_t : compute $(\hat{\mathbf{x}}_t)_{\hat{\mathcal{T}}_t} = ((\boldsymbol{\Phi}_t)_{\hat{\mathcal{T}}_t})^\dagger \mathbf{y}_t, (\hat{\mathbf{x}}_t)_{\hat{\mathcal{T}}_t^c} = 0$
2. Estimate $\boldsymbol{\ell}_t$: $\hat{\boldsymbol{\ell}}_t \leftarrow \mathbf{m}_t - \hat{\mathbf{x}}_t$.
3. If $t \bmod \alpha \neq 0$ then $\hat{\mathbf{P}}_{t,*} \leftarrow \hat{\mathbf{P}}_{t-1,*}, \hat{\mathbf{P}}_{t,\text{new}} \leftarrow \hat{\mathbf{P}}_{t-1,\text{new}}, \hat{\mathbf{P}}_t \leftarrow [\hat{\mathbf{P}}_{t,*} \hat{\mathbf{P}}_{t,\text{new}}]$
4. If $t \bmod \alpha = 0$ then detection or projection PCA

If phase = detect then

 - (a) Set $u = \frac{t}{\alpha}$ and compute $\mathcal{M}_u = (\mathbf{I} - \hat{\mathbf{P}}_{u\alpha-1,*} \hat{\mathbf{P}}_{u\alpha-1,*}') \left(\frac{1}{\alpha} \sum_{\tau=(u-1)\alpha+1}^{u\alpha} \hat{\boldsymbol{\ell}}_\tau \hat{\boldsymbol{\ell}}_\tau' \right) (\mathbf{I} - \hat{\mathbf{P}}_{u\alpha-1,*} \hat{\mathbf{P}}_{u\alpha-1,*}')$
 - (b) Update $\hat{\mathbf{P}}_{t,*} \leftarrow \hat{\mathbf{P}}_{t-1,*}, \hat{\mathbf{P}}_{t,\text{new}} \leftarrow \hat{\mathbf{P}}_{t-1,\text{new}}, \hat{\mathbf{P}}_t \leftarrow \hat{\mathbf{P}}_{t-1}$
 - (c) If $\lambda_{\max}(\mathcal{M}_u) \geq \text{thresh}$ then
 - i. phase \leftarrow ppca, increment $j \leftarrow j + 1$, reset $k \leftarrow 0$
 - ii. $\hat{u}_j \leftarrow u, \hat{t}_j = t$

Else (phase = ppca) then

 - (a) Set $u = \frac{t}{\alpha}$ and compute $\mathcal{M}_u = (\mathbf{I} - \hat{\mathbf{P}}_{u\alpha-1,*} \hat{\mathbf{P}}_{u\alpha-1,*}') \left(\frac{1}{\alpha} \sum_{\tau=(u-1)\alpha+1}^{u\alpha} \hat{\boldsymbol{\ell}}_\tau \hat{\boldsymbol{\ell}}_\tau' \right) (\mathbf{I} - \hat{\mathbf{P}}_{u\alpha-1,*} \hat{\mathbf{P}}_{u\alpha-1,*}')$
 - (b) Increment $k \leftarrow k + 1$
 - (c) Update $\hat{\mathbf{P}}_{t,\text{new}} \leftarrow \text{eigenvectors}(\mathcal{M}_u, \text{thresh}), \hat{\mathbf{P}}_{t,*} \leftarrow \hat{\mathbf{P}}_{t-1,*}$ and $\hat{\mathbf{P}}_t \leftarrow [\hat{\mathbf{P}}_{t,*} \hat{\mathbf{P}}_{t,\text{new}}]$.
Set $\hat{c}_{j,\text{new},k} \leftarrow \text{rank}(\hat{\mathbf{P}}_{t,\text{new}})$.
 - (d) If $k = K$
 - i. Set phase \leftarrow detect
 - ii. Update $\hat{\mathbf{P}}_{t,*} \leftarrow \hat{\mathbf{P}}_t$, and $\hat{\mathbf{P}}_{t,\text{new}} \leftarrow [\cdot]$

The function $\text{eigenvectors}(\mathcal{M}, \text{thresh})$ returns a basis matrix for the span of all eigenvectors whose eigenvalue is above thresh.

3.5 Proof of Theorem 3.2.15

We will prove Theorem 3.3.2 which will imply Theorem 3.2.15.

3.5.1 Definitions

Definition 3.5.1. *Define*

$$\mathcal{J}_u := [(u-1)\alpha + 1, u\alpha].$$

Also define u_j to be the u such that $t_j \in \mathcal{J}_u$. That is $u_j := \left\lceil \frac{t_j}{\alpha} \right\rceil$. For the purposes of describing events before the first subspace change, let $u_0 := 0$.

Also define $\hat{u}_j := \frac{t_j}{\alpha}$. Notice from the algorithm that this will be an integer, because detection only occurs when $t \bmod \alpha = 0$.

We will show that with high probability, $\hat{u}_j = u_j$ or $\hat{u}_j = u_j + 1$.

Recall that

1. $\mathbf{P}_{(j),*} := \mathbf{P}_{(j-1)} = \mathbf{P}_{t_{j-1}}$
2. $\mathbf{P}_{(j),\text{new}} := \mathbf{P}_{t_j,\text{new}}$.

Definition 3.5.2. *For $j = 1, 2, \dots, J$ and $k = 1, 2, \dots, K$ define*

1. $\hat{\mathbf{P}}_{(j),*} := \hat{\mathbf{P}}_{\hat{t}_{j-1} + K\alpha}$. This is the final estimate of $\mathbf{P}_{(j),*} = \mathbf{P}_{(j-1)}$.
2. $\hat{\mathbf{P}}_{(j),\text{new},0} := [\cdot]$ and $\hat{\mathbf{P}}_{(j),\text{new},k} := \hat{\mathbf{P}}_{\hat{t}_j + k\alpha,\text{new}}$ is the k^{th} estimate of $\mathbf{P}_{(j),\text{new}}$.

Notice from the algorithm that

1. $\hat{\mathbf{P}}_{t,*} = \hat{\mathbf{P}}_{(j),*}$ for all $t \in [\hat{t}_j, \hat{t}_j + K\alpha - 1]$
2. $\hat{\mathbf{P}}_{t,\text{new}} = \hat{\mathbf{P}}_{(j),\text{new},k}$ for all $t \in \mathcal{J}_{\hat{u}_j + (k+1)\alpha}$
3. At all times $\hat{\mathbf{P}}_t = [\hat{\mathbf{P}}_{t,*} \ \hat{\mathbf{P}}_{t,\text{new}}]$. Thus $\hat{\mathbf{P}}_t$ and $\hat{\mathbf{P}}_{t,\text{new}}$ update at every $t = \hat{t}_j + k\alpha$, $k = 1, 2, \dots, K$, $j = 1, 2, \dots, J$ while $\hat{\mathbf{P}}_{t,*}$ updates at every $t = \hat{t}_{j-1} + K\alpha$, $j = 2, \dots, J$.
4. $\hat{\mathbf{P}}_{t-1,*} \perp \hat{\mathbf{P}}_{t,\text{new}}$
5. $\Phi_t = (\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}' - \hat{\mathbf{P}}_{(j),\text{new},k} \hat{\mathbf{P}}_{(j),\text{new},k}')$ when $t \in [\hat{u}_j + k\alpha + 1, \hat{u}_j + (k+1)\alpha]$

Definition 3.5.3. Recall that for basis matrices \mathbf{P} and \mathbf{Q} , $\text{dif}(\mathbf{P}, \mathbf{Q}) := \|(I - \mathbf{P}\mathbf{P}')\mathbf{Q}\|_2$.

Define

1. $\zeta_{j,*} := \text{dif}(\hat{\mathbf{P}}_{(j),*}, \mathbf{P}_{(j),*})$
2. $\zeta_{j,\text{new},k} := \text{dif}([\hat{\mathbf{P}}_{(j),*} \ \hat{\mathbf{P}}_{(j),\text{new},k}], \mathbf{P}_{(j),\text{new}})$

Recall $\text{SE}_t = \text{dif}(\hat{\mathbf{P}}_t, \mathbf{P}_t)$, and notice that for $t \in \mathcal{J}_{\hat{u}+k}$, $\text{SE}_t \leq \zeta_{j,*} + \zeta_{j,\text{new},k-1}$.

Definition 3.5.4. Define

1. $\zeta_{j,*}^+ := (r_0 + (j-1)c_{\text{new}})\zeta$
2. $\zeta_{0,\text{new}}^+ := 1$, $\zeta_{k,\text{new}}^+ := \frac{b_{\mathbf{H},k}}{b_{\mathbf{A}} - b_{\mathbf{A},\perp} - b_{\mathbf{H},k}}$ (the right hand side depends on $\zeta_{k-1,\text{new}}^+$)

where $b_{\mathbf{A}}$, $b_{\mathbf{A},\perp}$, and $b_{\mathbf{H},k}$ are defined in (3.11), (3.12), and (3.13), respectively.

Definition 3.5.5. Define the random variable

$$X_u := \{\nu_1, \dots, \nu_{u\alpha}\}$$

If assuming Signal Model 3.2.17 (random support change), then both $\{\theta_t\}$ and $\{\mu_t\}$ for $t = 1, \dots, t_{\max}$ are also included in the definition of X_u for all u . Thus whenever conditioning on an X_u , they can be treated as deterministic.

Definition 3.5.6. For $j = 1, \dots, J$, $k = 1, \dots, K$, and for $a = u_j$ or $a = u_j + 1$, define the following events

- $\text{DET}_j^a := \{\hat{u}_j = a\}$
- $\text{PPCA}_{j,k}^a := \left\{ \hat{u}_j = a \text{ and } \text{rank}(\hat{\mathbf{P}}_{(j),\text{new},k}) = c_{j,\text{new}} \text{ and } \zeta_{j,\text{new},k} \leq \zeta_{k,\text{new}}^+ \right\}$
- $\text{NODETS}_j^a := \{\hat{u}_j = a \text{ and } \lambda_{\max}(\mathbf{M}_u) < \text{thresh for all } u \in [\hat{u}_j + K + 1, u_{j+1} - 1]\}$
- $\Gamma_{j,0}^a := \Gamma_{j-1,\text{end}} \cap \text{DET}_j^a$
- $\Gamma_{j,k}^a := \Gamma_{j,k-1}^a \cap \text{PPCA}_{j,k}^a$
- $\Gamma_{j,\text{end}} := \left(\Gamma_{j,K}^{u_j} \cap \text{NODETS}_j^{u_j} \right) \cup \left(\Gamma_{j,K}^{u_j+1} \cap \text{NODETS}_j^{u_j+1} \right)$

- $\Gamma_{0,\text{end}} := \{\zeta_{1,*} \leq r_0\zeta\} \cap \{\lambda_{\max}(\mathcal{M}_u) < \text{thresh for all } u \in [1, u_1 - 1]\}$

We misuse notation as follows. Suppose that a set Γ is a subset of all possible values that a r.v. X can take. For two r.v.s' $\{X, Y\}$, when we need to say “ $X \in \Gamma$ and Y can be anything” we will sometimes misuse notation and just say “ $\{X, Y\} \in \Gamma$.” For example, we sometimes say $X_{u_j} \in \Gamma_{j,\text{end}}$. This means $X_{u_j-1} \in \Gamma_{j,\text{end}}$ and $\boldsymbol{\nu}_t$ for $t \in \mathcal{J}_{u_j}$ are unconstrained.

Definition 3.5.7. Define \mathbf{e}_t to be the error made in estimating \mathbf{x}_t and ℓ_t . That is

$$\mathbf{e}_t := \hat{\mathbf{x}}_t - \mathbf{x}_t = \ell_t - \hat{\ell}_t$$

3.5.2 Main Lemmas

Fact 3.5.8. Under the assumptions of Theorem 3.2.15, $\frac{1}{\alpha} \leq (c_{\text{new}}\zeta)^2$. To see this, observe that the lower bound for α has $(c_{\text{new}}\zeta)^2$ in the denominator, and everything else in the expression is greater than or equal to 1. (Notice that $\frac{\gamma_{\text{new}}^2}{\lambda^-} \geq 1$)

Lemma 3.5.9. Under the conditions of Theorem 3.2.15,

$$\zeta_{k,\text{new}}^+ \leq 0.4^k + 0.84c_{\text{new}}\zeta$$

Proof. This claim follows by applying simple algebra on the definition and using the bounds assumed on α , ζ , and b in Theorem 3.2.15. In particular, we use the fact that $\frac{1}{\alpha} \leq (c_{\text{new}}\zeta)^2$ (Fact 3.5.8). Detailed proofs of similar results can be found in [12] and [13, Lemma 6.2]. Those proofs define the quantity $g := \frac{\lambda_{\text{new}}^+}{\lambda_{\text{new}}^-}$ and assume $g \leq \sqrt{2}$. Although we do not explicitly define g , our assumption that $\lambda_{\text{new}}^+ \leq \sqrt{2}\lambda^-$ implies that $\frac{\lambda_{\text{new}}^+}{\lambda_{\text{new}}^-} \leq \sqrt{2}$, because $\lambda_{\text{new}}^- \geq \lambda^-$.

For the purposes of review, we have included a proof of this lemma in Appendix 3.C. \square

Lemma 3.5.10 (Sparse Recovery Lemma (Chapter 2, Lemma 2.6.15)). Assume that all of the conditions of Theorem 3.2.15 hold. Recall that $\text{SE}_t = \text{dif}(\hat{\mathbf{P}}_t, \mathbf{P}_t)$.

1. Conditioned on $\Gamma_{j-1,\text{end}}$, for $t \in [t_j, (u_j + 1)\alpha]$

$$(a) \ \phi_t := \|[(\Phi_t)_{\mathcal{T}_t}'(\Phi_t)_{\mathcal{T}_t}]^{-1}\|_2 \leq \phi^+ := 1.2.$$

(b) the support of \mathbf{x}_t is recovered exactly i.e. $\hat{\mathcal{T}}_t = \mathcal{T}_t$ and \mathbf{e}_t satisfies:

$$\mathbf{e}_t := \hat{\mathbf{x}}_t - \mathbf{x}_t = \ell_t - \hat{\ell}_t = \mathbf{I}_{\mathcal{T}_t}[(\Phi_t)_{\mathcal{T}_t}'(\Phi_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \Phi_t \ell_t. \quad (3.9)$$

(c) Furthermore,

$$\text{SE}_t \leq 1, \text{ and}$$

$$\|\mathbf{e}_t\|_2 \leq \phi^+(\zeta_{j,*}^+ \sqrt{r}\gamma + \sqrt{c_{\text{new}}}\gamma_{\text{new}}) \leq 1.2 \left(\sqrt{\zeta} + \sqrt{c_{\text{new}}}\gamma_{\text{new}} \right)$$

2. For $k = 1, 2, \dots, K$ and $\hat{u}_j = u_j$ or $\hat{u}_j = u_j + 1$, conditioned on $\Gamma_{j,k-1}^{\hat{u}_j}$, for $t \in [(\hat{u}_j + k - 1)\alpha + 1, (\hat{u}_j + k)\alpha]$, the first two conclusions above hold. That is, $\phi_t \leq \phi^+$ and \mathbf{e}_t satisfies (3.9). Furthermore,

$$\text{SE}_t \leq \zeta_{j,*}^+ + \zeta_{j,\text{new},k-1}^+, \text{ and}$$

$$\|\mathbf{e}_t\|_2 \leq \phi^+(\zeta_{j,*}^+ \sqrt{r}\gamma + \zeta_{j,\text{new},k-1}^+ \sqrt{c_{\text{new}}}\gamma_{\text{new}}) \leq 1.2 \left(1.84\sqrt{\zeta} + (0.4)^{k-1} \sqrt{c_{\text{new}}}\gamma_{\text{new}} \right)$$

3. For $\hat{u}_j = u_j$ or $\hat{u}_j = u_j + 1$, conditioned on $\Gamma_{j,K}^{\hat{u}_j}$, for $t \in [(\hat{u}_j + K)\alpha + 1, t_{j+1} - 1]$, the first two conclusions above hold ($\phi_t \leq \phi^+$ and \mathbf{e}_t satisfies (3.9)). Furthermore,

$$\text{SE}_t \leq \zeta_{j+1,*}^+, \text{ and}$$

$$\|\mathbf{e}_t\|_2 \leq \phi^+ \zeta_{j+1,*}^+ \sqrt{r}\gamma \leq 1.2\sqrt{\zeta}$$

The proof in Chapter 2 uses the facts that under the various conditionings, $\text{rank}(\hat{\mathbf{P}}_{(j),k-1,\text{new}}) = c_{j,\text{new}}$, $\zeta_{j,*} \leq \zeta_{j,*}^+$, and $\zeta_{j,\text{new},k-1} \leq \zeta_{k-1,\text{new}}^+$.

Lemma 3.5.11.

1. The event $\Gamma_{j,K}^{\hat{u}_j}$ and so also the event $\Gamma_{j,\text{end}}$ imply that $\zeta_{j+1,*} \leq \zeta_{j+1,*}^+$.
2. $\mathbb{P}(\text{NODETS}_j^a \mid \Gamma_{j,K}^a) = 1$ for $a = u_j$ or $a = u_j + 1$.

Lemma 3.5.12. For $j = 1, \dots, J$,

$$\mathbb{P}(\text{DET}^{u_j+1} \mid \Gamma_{j-1,\text{end}}, \overline{\text{DET}^{u_j}}) \geq p_{\text{det},1} := 1 - p_{\mathbf{A}} - p_{\mathbf{H}}.$$

The definitions of $p_{\mathbf{A}}$ and $p_{\mathbf{H}}$ can be found in Lemmas 3.5.20 and 3.5.22 respectively.

Fact 3.5.13. *Observe that $\Gamma_{j,0}^a$ for $a = u_j$ or $a = u_j + 1$ implies that $u_j \leq \hat{u}_j \leq u_j + 1$. Thus, $t_j \leq \hat{t}_j \leq t_j + 2\alpha$. So with the model assumption that $d \geq (K + 2)\alpha$, we have that $\mathcal{J}_{\hat{u}_j+k} \subseteq [t_j, t_j + d]$ for $k = 1, 2, \dots, K$. This fact is needed so that the tighter bounds on $\mathbf{a}_{t,\text{new}}$ hold from condition 3b in the Theorem.*

Lemma 3.5.14.

$$\mathbb{P}(\Gamma_{j,k}^a \mid \Gamma_{j,k-1}^a) = \mathbb{P}(\text{PPCA}_{j,k}^a \mid \Gamma_{j,k-1}^a) \geq p_{\text{ppca}} := 1 - p_{\mathbf{A}} - p_{\mathbf{A},\perp} - p_{\mathcal{H}}$$

for $a = u_j$ or $a = u_j + 1$. where $p_{\text{ppca}} := 1 - p_{\mathbf{A}} - p_{\mathbf{A},\perp} - p_{\mathcal{H}}$. The definitions of $p_{\mathbf{A}}$, $p_{\mathbf{A},\perp}$, and $p_{\mathcal{H}}$ can be found in Lemmas 3.5.20, 3.5.21, and 3.5.22 respectively.

The above lemma says that whether the new directions are detected at u_j or $u_j + 1$, conditioned on $k - 1$ previous successful projection PCA steps, the probability of a successful k^{th} projection PCA step is lower bounded by p_{ppca} .

Corollary 3.5.15. *Let*

$$p_{\text{det},0} := \mathbb{P}(\text{DET}^{u_j} \mid \Gamma_{j-1,\text{end}})$$

and therefore, $1 - p_{\text{det},0} = \mathbb{P}(\overline{\text{DET}^{u_j}} \mid \Gamma_{j-1,\text{end}})$.

From the above lemmas, we get that

$$\begin{aligned} \mathbb{P}(\Gamma_{j,\text{end}} \mid \Gamma_{j-1,\text{end}}) &= \mathbb{P}\left((\text{DET}^{u_j} \cap \text{PPCA}_{j,1}^{u_j} \cap \dots \cap \text{PPCA}_{j,K}^{u_j} \cap \text{NODETS}_j^{u_j}) \cup \right. \\ &\quad \left. (\overline{\text{DET}^{u_j}} \cap \text{DET}^{u_j+1} \cap \text{PPCA}_{j,k}^{u_j+1} \cap \dots \cap \text{PPCA}_{j,K}^{u_j+1} \cap \text{NODETS}_j^{u_j+1}) \mid \Gamma_{j-1,\text{end}}\right) \\ &= \mathbb{P}\left(\text{DET}^{u_j} \cap \text{PPCA}_{j,1}^{u_j} \cap \dots \cap \text{PPCA}_{j,K}^{u_j} \mid \Gamma_{j-1,\text{end}}\right) \\ &\quad + \mathbb{P}\left(\overline{\text{DET}^{u_j}} \cap \text{DET}^{u_j+1} \cap \text{PPCA}_{j,k}^{u_j+1} \cap \dots \cap \text{PPCA}_{j,K}^{u_j+1} \mid \Gamma_{j-1,\text{end}}\right) \\ &\geq p_{\text{det},0} \cdot (p_{\text{ppca}})^K + (1 - p_{\text{det},0}) \cdot p_{\text{det},1} \cdot (p_{\text{ppca}})^K \\ &\geq p_{\text{det},0} \cdot p_{\text{det},1} \cdot (p_{\text{ppca}})^K + (1 - p_{\text{det},0}) \cdot p_{\text{det},1} \cdot (p_{\text{ppca}})^K \\ &= p_{\text{det},1} \cdot (p_{\text{ppca}})^K. \end{aligned}$$

Proof of Theorem 3.3.2. Theorem 3.3.2 follows from Corollary 3.5.15 and the assumed lower bound on α . Notice that by Lemma 3.5.9, the choice of K , and Lemma 3.5.10, the event $\Gamma_{J,\text{end}}$ will imply all conclusions of the theorem.

By the first assumption of Theorem 3.2.15 and the argument used to prove Lemma 3.5.11, we get that $\mathbb{P}(\Gamma_{0,\text{end}}) = 1$.

Next, we have by the chain rule,

$$\mathbb{P}(\Gamma_{J,\text{end}}) = \prod_{j=1}^J \mathbb{P}(\Gamma_{j,\text{end}} \mid \Gamma_{j-1,\text{end}}, \Gamma_{j-2,\text{end}}, \dots, \Gamma_{1,\text{end}}, \Gamma_{0,\text{end}}).$$

Because $\Gamma_{j-1,\text{end}} \subseteq \Gamma_{j-2,\text{end}} \subseteq \dots \subseteq \Gamma_{1,\text{end}} \subseteq \Gamma_{0,\text{end}}$, we get

$$\begin{aligned} \mathbb{P}(\Gamma_{J,\text{end}}) &= \prod_{j=1}^J \mathbb{P}(\Gamma_{j,\text{end}} \mid \Gamma_{j-1,\text{end}}) \\ &\geq \prod_{j=1}^J p_{\text{det},1} \cdot (p_{\text{ppca}})^K \\ &= (p_{\text{det},1} \cdot (p_{\text{ppca}})^K)^J \\ &\geq 1 - n^{-10} \end{aligned}$$

The last line is by the lower bound on α assumed in Theorem 3.2.15. \(\square\)

3.5.3 Key Lemmas for Proving of Lemmas 3.5.11, 3.5.12, and 3.5.14

Before proving the lemmas from the preceding section, we introduce several lemmas which will be used in the proofs. Their statement requires the following definition.

Definition 3.5.16. Recall from Algorithm 3 that $\mathcal{M}_u = \frac{1}{\alpha} \sum_{t \in \mathcal{J}_u} (\mathbf{I} - \hat{\mathbf{P}}_{u\alpha-1,*} \hat{\mathbf{P}}_{u\alpha-1,*}') \hat{\ell}_t \hat{\ell}_t' (\mathbf{I} - \hat{\mathbf{P}}_{u\alpha-1,*} \hat{\mathbf{P}}_{u\alpha-1,*}')$. Using the definition of $\mathbf{P}_{(j),*}$, for $u = u_j + 1$ or $u = \hat{u}_j + k$,

$$\mathcal{M}_u = (\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}') \left(\frac{1}{\alpha} \sum_{t \in \mathcal{J}_u} \hat{\ell}_t \hat{\ell}_t' \right) (\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}').$$

Define

1. Let $\mathbf{D}_{j,\text{new}} := (\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}') \mathbf{P}_{(j),\text{new}} \stackrel{QR}{=} \mathbf{E}_{j,\text{new}} \mathbf{R}_{j,\text{new}}$ denote its reduced QR decomposition, i.e. let $\mathbf{E}_{j,\text{new}}$ be a basis matrix for $\text{range}(\mathbf{D}_{j,\text{new}})$ and let $\mathbf{R}_{j,\text{new}} = \mathbf{E}_{j,\text{new}}' \mathbf{D}_{j,\text{new}}$.
2. Let $\mathbf{E}_{j,\text{new},\perp}$ be a basis matrix for the orthogonal complement of $\text{range}(\mathbf{E}_{j,\text{new}})$. To be precise, $\mathbf{E}_{j,\text{new},\perp}$ is an $n \times (n - r_j)$ basis matrix so that $[\mathbf{E}_{j,\text{new}} \ \mathbf{E}_{j,\text{new},\perp}]$ is unitary.

3. For $u = u_j + 1$ and $u = \hat{u}_j + k$ for $k = 1, \dots, K$, define \mathbf{A}_u , $\mathbf{A}_{u,\perp}$, \mathcal{A}_u and \mathcal{H}_u as

$$\begin{aligned}\mathbf{A}_u &:= \frac{1}{\alpha} \sum_{t \in \mathcal{J}_u} \mathbf{E}_{j,\text{new}}' (I - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}') \ell_t \ell_t' (I - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}') \mathbf{E}_{j,\text{new}} \\ \mathbf{A}_{u,\perp} &:= \frac{1}{\alpha} \sum_{t \in \mathcal{J}_u} \mathbf{E}_{j,\text{new},\perp}' (I - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}') \ell_t \ell_t' (I - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}') \mathbf{E}_{j,\text{new},\perp}\end{aligned}$$

and let

$$\mathcal{A}_u := \begin{bmatrix} \mathbf{E}_{j,\text{new}} & \mathbf{E}_{j,\text{new},\perp} \end{bmatrix} \begin{bmatrix} \mathbf{A}_u & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_{u,\perp} \end{bmatrix} \begin{bmatrix} \mathbf{E}_{j,\text{new}}' \\ \mathbf{E}_{j,\text{new},\perp}' \end{bmatrix}$$

and

$$\mathcal{H}_u := \mathcal{M}_u - \mathcal{A}_u$$

Recall that for the above values of u_j ,

$$\mathcal{M}_u = (I - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}') \left(\frac{1}{\alpha} \sum_{t \in \mathcal{J}_u} \hat{\ell}_t \hat{\ell}_t' \right) (I - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}').$$

When $u = \hat{u}_j + k$ for a $k \leq K$, \mathcal{M}_u is the matrix whose eigenvectors with eigenvalue above thresh form $\hat{\mathbf{P}}_{(j),\text{new},k}$ (see step 4c of Algorithm 3). In this case, \mathcal{M}_u has eigendecomposition

$$\mathcal{M}_u \stackrel{\text{EVD}}{=} \begin{bmatrix} \hat{\mathbf{P}}_{(j),\text{new},k} & \hat{\mathbf{P}}_{(j),k,\text{new},\perp} \end{bmatrix} \begin{bmatrix} \hat{\Lambda}_u & \mathbf{0} \\ \mathbf{0} & \hat{\Lambda}_{u,\perp} \end{bmatrix} \begin{bmatrix} \hat{\mathbf{P}}_{(j),\text{new},k}' \\ \hat{\mathbf{P}}_{(j),k,\text{new},\perp}' \end{bmatrix}.$$

Note: we use \mathbf{A}_u , $\mathbf{A}_{u,\perp}$, and \mathcal{H}_u for $u = u_j + 1$ or for $u = \hat{u}_j + 1, \dots, \hat{u}_j + K$. With the appropriate conditioning, all these u 's lie in the interval $[u_j + 1, u_j + d - 1]$ and from the assumptions in the Theorem, in this interval $\mathbf{a}_{t,\text{new}}$ is “small”.

The following lemma follows from the $\sin \theta$ Theorem [26] and Weyl's theorem. It is taken from [12].

Lemma 3.5.17 ([12]). *At $u = \hat{u}_j + k$, if $\text{rank}(\hat{\mathbf{P}}_{(j),\text{new},k}) = c_{j,\text{new}}$, and if $\lambda_{\min}(\mathbf{A}_u) - \|\mathbf{A}_{u,\perp}\|_2 - \|\mathcal{H}_u\|_2 > 0$, then*

$$\zeta_{j,\text{new},k} \leq \frac{\|\mathcal{H}_u\|_2}{\lambda_{\min}(\mathbf{A}_u) - \|\mathbf{A}_{u,\perp}\|_2 - \|\mathcal{H}_u\|_2} \quad (3.10)$$

The next three lemmas each assert a high probability bound for one of the terms in (3.10).

Definition 3.5.18. For ease of notation, define the function

$$F(\alpha, \epsilon, b_1) := \exp\left(\frac{-\alpha\epsilon^2}{8(b_1)^2}\right)$$

In the following lemmas, let

$$\epsilon = 0.01c_{\text{new}}\zeta\lambda^-.$$

Remark 3.5.19. In the next three lemmas we use quantities $b_{\mathbf{A}}$, $b_{\mathbf{A},\perp}$, and $b_{\mathbf{H},k}$. These are not to be confused with the autoregression (AR) coefficient b used for the AR model on \mathbf{a}_t (see Signal Model 3.2.1).

Let

$$p_{\mathbf{A}} := cF(\alpha, \epsilon, c_{\text{new}}\gamma_{\text{new}}^2) + 3(r + c_{\text{new}})F(\alpha, \epsilon, 2\sqrt{c_{\text{new}}r}\gamma_{\text{new}}) + (r + c_{\text{new}})F(\alpha, \epsilon, 4\sqrt{c_{\text{new}}r}\gamma_{\text{new}})$$

and

$$b_{\mathbf{A}} := (1 - (\zeta_{j,*}^+)^2)(1 - b^2)\lambda_{\text{new}}^- - \epsilon - 2\zeta_{j,*}^+ \left(\frac{b^2}{\alpha(1 - b^2)} \sqrt{c_{\text{new}}r}\gamma_{\text{new}}\gamma + 4\epsilon \right). \quad (3.11)$$

Lemma 3.5.20. For $k = 1, \dots, K$,

$$\mathbb{P}(\lambda_{\min}(\mathbf{A}_{\hat{u}_j+k}) \geq b_{\mathbf{A}} \mid X_{\hat{u}_j+k-1}) \geq 1 - p_{\mathbf{A}}$$

for all $X_{\hat{u}_j+k-1} \in \Gamma_{j,k-1}^{\hat{u}_j}$ with $\hat{u}_j = u_j$ or $\hat{u}_j = u_j + 1$. Also,

$$\mathbb{P}(\lambda_{\min}(\mathbf{A}_{u_j+1}) \geq b_{\mathbf{A}} \mid X_{u_j}) \geq 1 - p_{\mathbf{A}}$$

for all $X_{u_j} \in \Gamma_{j-1,\text{end}}$.

Let $p_{\mathbf{A},\perp} := rF(\alpha, \epsilon, (\zeta_{j,*}^+)^2 r \gamma^2)$ and

$$b_{\mathbf{A},\perp} := (\zeta_{j,*}^+)^2 (b^2 r \gamma^2 + (1 - b^2)\lambda^+ + \epsilon). \quad (3.12)$$

Lemma 3.5.21. For $k = 1, \dots, K$,

$$\mathbb{P}(\lambda_{\max}(\mathbf{A}_{\hat{u}_j+k,\perp}) \leq b_{\mathbf{A},\perp} \mid X_{\hat{u}_j+k-1}) \geq 1 - p_{\mathbf{A},\perp}$$

for all $X_{\hat{u}_j+k-1} \in \Gamma_{j,k-1}^{\hat{u}_j}$ with $\hat{u}_j = u_j$ or $\hat{u}_j = u_j + 1$.

The same bound holds for $\|\mathbf{A}_{u_j+1,\perp}\|_2$ when we condition on $X_{u_j} \in \Gamma_{j-1,\text{end}}$.

Let $p_{\mathcal{H}} := nF \left(\alpha, \epsilon, \left[\phi^+ \left(\zeta_{j,*}^+ \sqrt{r}\gamma + \sqrt{c_{\text{new}}}\gamma_{\text{new}} \right) \right]^2 \right) + rF \left(\alpha, \epsilon, r\gamma^2 \right) + c_{\text{new}}F \left(\alpha, \epsilon, c_{\text{new}}\gamma_{\text{new}}^2 \right) + 3(r + c_{\text{new}})F \left(\alpha, \epsilon, 2\sqrt{c_{\text{new}}r}\gamma\gamma_{\text{new}} \right) + (r + c_{\text{new}})F \left(\alpha, \epsilon, 4\sqrt{c_{\text{new}}r}\gamma\gamma_{\text{new}} \right) + c_{\text{new}}F \left(\alpha, \epsilon, r\gamma^2 \right) + 3(r + c_{\text{new}})F \left(\alpha, \epsilon, 2\sqrt{c_{\text{new}}r}\gamma\gamma_{\text{new}} \right) + (r + c_{\text{new}})F \left(\alpha, \epsilon, 4\sqrt{c_{\text{new}}r}\gamma\gamma_{\text{new}} \right)$. Define $\kappa_s^+ := .0215$. Also let

$$b_{\mathcal{H},k} := b_{2,k} + 2b_{4,k} + 2b_6, \quad (3.13)$$

where

$$b_{2,k} := \begin{cases} \begin{aligned} &\rho^2 h^+ (\phi^+)^2 (\zeta_{j,*}^+)^2 (b^2 r \gamma^2 + (1 - b^2) \lambda^+) + \\ &\rho^2 h^+ (\phi^+)^2 (\kappa_s^+)^2 (b^2 c_{\text{new}} \gamma_{\text{new}}^2 + (1 - b^2) \lambda_{\text{new}}^+) + \\ &2 \cdot \rho^2 h^+ (\phi^+)^2 \kappa_s^+ \zeta_{j,*}^+ b^2 \sqrt{r c_{\text{new}}} \gamma \gamma_{\text{new}} \end{aligned} & k = 1 \\ \begin{aligned} &\rho^2 h^+ (\phi^+)^2 (\zeta_{j,*}^+)^2 (b^2 r \gamma^2 + (1 - b^2) \lambda^+) + \\ &\rho^2 h^+ (\phi^+)^2 (\zeta_{j,\text{new},k-1}^+)^2 (b^2 c_{\text{new}} \gamma_{\text{new}}^2 + (1 - b^2) \lambda_{\text{new}}^+) + \\ &2 \cdot \rho^2 h^+ (\phi^+)^2 \zeta_{j,*}^+ \zeta_{j,\text{new},k-1}^+ b^2 \sqrt{r c_{\text{new}}} \gamma \gamma_{\text{new}} \end{aligned} & k \geq 2 \end{cases}$$

$$b_{4,k} := \begin{cases} \begin{aligned} &\left[(\zeta_{j,*}^+)^2 (b^2 r \gamma^2 + (1 - b^2) \lambda^+ + \epsilon) + \right. \\ &\kappa_s^+ (b^2 c_{\text{new}} \gamma_{\text{new}}^2 + (1 - b^2) \lambda_{\text{new}}^+ + \epsilon) + \\ &\zeta_{j,*}^+ \kappa_s^+ \left(\frac{b^2}{\alpha(1-b^2)} \sqrt{c_{\text{new}}r} \gamma_{\text{new}} \gamma + 4\epsilon \right) + \\ &\left. \zeta_{j,*}^+ \left(\frac{b^2}{\alpha(1-b^2)} \sqrt{c_{\text{new}}r} \gamma_{\text{new}} \gamma + 4\epsilon \right) \right] \phi^+ \end{aligned} & k = 1 \\ \begin{aligned} &\left[(\zeta_{j,*}^+)^2 (b^2 r \gamma^2 + (1 - b^2) \lambda^+ + \epsilon) \right. \\ &+ \zeta_{j,\text{new},k-1}^+ (b^2 c_{\text{new}} \gamma_{\text{new}}^2 + (1 - b^2) \lambda_{\text{new}}^+ + \epsilon) \\ &\zeta_{j,*}^+ \zeta_{j,\text{new},k-1}^+ \left(\frac{b^2}{\alpha(1-b^2)} \sqrt{c_{\text{new}}r} \gamma_{\text{new}} \gamma + 4\epsilon \right) \\ &\left. + \zeta_{j,*}^+ \left(\frac{b^2}{\alpha(1-b^2)} \sqrt{c_{\text{new}}r} \gamma_{\text{new}} \gamma + 4\epsilon \right) \right] \left(\sqrt{\rho^2 h^+} \phi^+ \right) \end{aligned} & k \geq 2 \end{cases}$$

and

$$b_6 := (\zeta_{j,*}^+)^2 (b^2 r \gamma^2 + (1 - b^2) \lambda^+ + \epsilon) + \zeta_{j,*}^+ \left(\frac{b^2}{\alpha(1-b^2)} \sqrt{c_{\text{new}}r} \gamma_{\text{new}} \gamma + 4\epsilon \right).$$

Lemma 3.5.22. For $k = 1, \dots, K$,

$$\mathbb{P}(\|\mathcal{H}_{\hat{u}_j+k}\|_2 \leq b_{\mathcal{H},k} + \epsilon \mid X_{\hat{u}_j+k-1}) \geq 1 - p_{\mathcal{H}} \quad (3.14)$$

for all $X_{\hat{u}_j+k-1} \in \Gamma_{j,k-1}^{\hat{u}_j}$ with $\hat{u}_j = u_j$ or $\hat{u}_j = u_j + 1$

The same bound (with $k = 1$) holds for $\|\mathcal{H}_{u_j+1}\|_2$ when we condition on $X_{u_j} \in \Gamma_{j-1,\text{end}}$.

The above lemmas are proved in the next section (Section 3.6). The proofs use Fact 3.5.13.

3.5.4 Proofs of Lemmas 3.5.11, 3.5.12, and 3.5.14

Proof of Lemma 3.5.11. Recall that $\Gamma_{j,\text{end}} := \left(\Gamma_{j,K}^{u_j} \cap \text{NODETS}_j^{u_j}\right) \cup \left(\Gamma_{j,K}^{u_j+1} \cap \text{NODETS}_j^{u_j+1}\right)$.

1. By the definition of $\Gamma_{j,K}^{\hat{u}_j}$ for $\hat{u} = u_j$ or $\hat{u} = u_j + 1$, $\zeta_{j,*} \leq \zeta_{j,*}^+ = (r_0 + (j-1)c_{\text{new}})\zeta$ and $\zeta_{j,K} \leq \zeta_{j,\text{new},K}^+$. Lemma 3.5.9 and the choice of K imply that $\zeta_{j,\text{new},K}^+ \leq c_{\text{new}}\zeta$. Thus, $\zeta_{j+1,*} \leq \zeta_{j+1,*}^+ = (r_0 + jc_{\text{new}})\zeta$. This holds for $\hat{u}_j = u_j$ or $\hat{u}_j = u_j + 1$.
2. $\mathbb{P}(\text{NODETS}_j^{\hat{u}_j} \mid \Gamma_{j,K}^{\hat{u}_j}) = \mathbb{P}(\lambda_{\max}(\mathcal{M}_u) < \text{thresh for all } u \in [u_j + K + 1, u_{j+1} - 1] \mid \Gamma_{j,K}^{\hat{u}_j})$
for $\hat{u}_j = u_j$ or $\hat{u}_j = u_j + 1$.

As shown in 1), $\Gamma_{j,K}^{\hat{u}_j}$ implies that $\text{dif}(\hat{P}_{(j+1),*}, \mathbf{P}_{(j+1),*}) \leq \zeta_{j+1,*}^+ = (r_0 + jc_{\text{new}})\zeta$. Notice that for $u \in [\hat{u}_j + K + 1, u_{j+1} - 1]$, $\hat{P}_{u\alpha-1,*} = \hat{P}_{(j+1),*}$, $\mathbf{P}_{(j+1),*} = \mathbf{P}_{(j)}$, and $\ell_t = \mathbf{P}_{(j)}\mathbf{a}_t$. Therefore,

$$\begin{aligned} \lambda_{\max}(\mathcal{M}_u) &= \lambda_{\max}\left(\frac{1}{\alpha} \sum_{t \in \mathcal{J}_u} (\mathbf{I} - \hat{P}_{u\alpha-1,*} \hat{P}_{u\alpha-1,*}') \hat{\ell}_t \hat{\ell}_t' (\mathbf{I} - \hat{P}_{u\alpha-1,*} \hat{P}_{u\alpha-1,*}')\right) \\ &= \lambda_{\max}\left(\frac{1}{\alpha} \sum_{t \in \mathcal{J}_u} (\mathbf{I} - \hat{P}_{(j+1),*} \hat{P}_{(j+1),*}') (\mathbf{P}_{(j)} \mathbf{a}_t - \mathbf{e}_t) \right. \\ &\quad \left. (\mathbf{P}_{(j)} \mathbf{a}_t - \mathbf{e}_t)' (\mathbf{I} - \hat{P}_{(j+1),*} \hat{P}_{(j+1),*}')\right) \\ &\leq (\zeta_{j+1,*}^+)^2 r \gamma^2 + 2\phi^+ (\zeta_{j+1,*}^+)^2 r \gamma^2 + (\phi^+)^2 (\zeta_{j+1,*}^+)^2 r \gamma^2 \\ &\leq 4(\phi^+)^2 \zeta \lambda^- \\ &\leq \frac{\lambda^-}{2}. \end{aligned}$$

The bounds on \mathbf{e}_t come from Lemma 3.5.10, and the penultimate line uses the bounds assumed on ζ in Theorem 3.2.15.

□

Fact 3.5.23. *Using the bounds on ζ , $b \leq 0.1$, $\lambda_{\text{new}}^+ \leq \sqrt{2}\lambda^-$, $c_{\text{new}}\eta_{\text{new}}\lambda_{\text{new}}^+ \leq 50\lambda^-$ from the Theorem and Fact 3.5.8 (which gives $\frac{1}{\alpha} \leq (c_{\text{new}}\zeta)^2$), we get:*

$$\begin{aligned} b_{\mathbf{A}} &\geq 0.94\lambda_{\text{new}}^- \geq 0.94\lambda^- \\ b_{\mathbf{A},\perp} &\leq 0.011\lambda^- \\ b_{\mathbf{H},k} &\leq 0.24\lambda^-. \end{aligned}$$

Thus, $b_{\mathbf{A}} - b_{\mathbf{H},k} > 0.5\lambda^- = \text{thresh}$ and $b_{\mathbf{A},\perp} + b_{\mathbf{H},k} < \text{thresh}$.

Proof of Lemma 3.5.12. We will prove that $\mathbb{P}(\text{DET}^{u_j+1} \mid X_{u_j}) > p_{\text{det},1}$ for all $X_{u_j} \in \Gamma_{j-1,\text{end}}$. In particular, this will imply that $\mathbb{P}(\text{DET}^{u_j+1} \mid X_{u_j}) > p_{\text{det},1}$ for all $X_{u_j} \in \Gamma_{j-1,\text{end}} \cap \overline{\text{DET}^{u_j}}$ and so we can conclude that $\mathbb{P}(\text{DET}^{u_j+1} \mid \Gamma_{j-1,\text{end}}, \overline{\text{DET}^{u_j}}) > p_{\text{det},1}$.

Observe that

$$\mathbb{P}(\text{DET}^{u_j+1} \mid X_{u_j}) = \mathbb{P}(\lambda_{\max}(\mathbf{M}_{u_j+1}) > \text{thresh} \mid X_{u_j})$$

By Weyl's Theorem

$$\begin{aligned} \lambda_{\max}(\mathbf{M}_{u_j+1}) &\geq \lambda_{\max}(\mathbf{A}_{u_j+1}) + \lambda_{\min}(\mathbf{H}_{u_j+1}) \\ &\geq \lambda_{\max}(\mathbf{A}_{u_j+1}) - \|\mathbf{H}_{u_j+1}\|_2 \\ &\geq \lambda_{\min}(\mathbf{A}_{u_j+1}) - \|\mathbf{H}_{u_j+1}\|_2 \\ &\geq b_{\mathbf{A}} - b_{\mathbf{H},1} \\ &\geq \frac{\lambda^-}{2} \end{aligned}$$

When $X_{u_j} \in \Gamma_{j-1,\text{end}}$, Lemmas 3.5.20 and 3.5.22 give high probability bounds on $\lambda_{\min}(\mathbf{A}_{u_j+1})$ and $\|\mathbf{H}_{u_j+1}\|_2$ respectively. So the above inequality holds with probability greater than or equal to $1 - p_{\mathbf{A}} - p_{\mathbf{H}} = p_{\text{det},1}$. □

Proof of Lemma 3.5.14. To prove this Lemma we need to show two things. First, conditioned on $\Gamma_{j,k-1}^{\hat{u}_j}$, the k^{th} estimate of the number of new directions is correct. That is: $\hat{c}_{j,\text{new},k} = c_{j,\text{new}}$. Second, we must show $\zeta_{j,\text{new},k} \leq \zeta_{j,\text{new},k}^+$, again conditioned on $\Gamma_{j,k-1}^{\hat{u}_j}$.

Notice that $\hat{c}_{j,\text{new},k} = \text{rank}(\hat{\mathbf{P}}_{(j),\text{new},k})$. To show that $\text{rank}(\hat{\mathbf{P}}_{(j),\text{new},k}) = c_{j,\text{new}}$, we need to show that for $u = \hat{u}_j + k$, $k = 1, \dots, K$, $\lambda_{c_{j,\text{new}}}(\mathbf{M}_u) > \text{thresh}$ and $\lambda_{c_{j,\text{new}}+1}(\mathbf{M}_u) < \text{thresh}$. To do this we proceed similarly to above. Observe that, $\mathbf{M}_u = \mathbf{A}_u + \mathbf{H}_u$. By Weyl's Theorem

$$\begin{aligned}\lambda_{c_{j,\text{new}}}(\mathbf{M}_u) &\geq \lambda_{c_{j,\text{new}}}(\mathbf{A}_u) + \lambda_{\min}(\mathbf{H}_u) \\ &\geq \lambda_{c_{j,\text{new}}}(\mathbf{A}_u) - \|\mathbf{H}_u\|_2 \\ &= \lambda_{\min}(\mathbf{A}_u) - \|\mathbf{H}_u\|_2.\end{aligned}$$

The equality is because \mathbf{A}_u is of size $c_{j,\text{new}} \times c_{j,\text{new}}$ and $\lambda_{\min}(\mathbf{A}_u) > \lambda_{\max}(\mathbf{A}_{u,\perp})$. Similarly,

$$\begin{aligned}\lambda_{c_{j,\text{new}}+1}(\mathbf{M}_u) &\leq \lambda_{c_{j,\text{new}}+1}(\mathbf{A}_u) + \lambda_{\max}(\mathbf{H}_u) \\ &\leq \lambda_{c_{j,\text{new}}+1}(\mathbf{A}_u) + \|\mathbf{H}_u\|_2 \\ &= \lambda_{\max}(\mathbf{A}_{u,\perp}) + \|\mathbf{H}_u\|_2.\end{aligned}$$

Using Lemmas 3.5.20, 3.5.21, and 3.5.22 and Fact 3.5.23, we can conclude that with probability greater than $1 - p_{\text{ppca}}$, $\lambda_{c_{j,\text{new}}}(\mathbf{M}_u) > b_{\mathbf{A}} - b_{\mathbf{H},k} \geq \lambda^-/2 = \text{thresh}$ and $\lambda_{c_{j,\text{new}}+1}(\mathbf{M}_u) < b_{\mathbf{A},\perp} + b_{\mathbf{H},k} \leq \lambda^-/2 = \text{thresh}$. Therefore $\text{rank}(\hat{\mathbf{P}}_{(j),\text{new},k}) = c_{j,\text{new}}$ with probability greater than $1 - p_{\text{ppca}}$.

To show that $\zeta_{j,\text{new},k} \leq \zeta_{k,\text{new}}^+$, we also use Lemmas 3.5.20, 3.5.21, and 3.5.22. Using $\text{rank}(\hat{\mathbf{P}}_{(j),\text{new},k}) = c_{j,\text{new}}$ and applying Lemma 3.5.17 with these bounds gives the desired result.

□

3.6 Proofs of Lemmas 3.5.20, 3.5.21, and 3.5.22

3.6.1 Key Lemmas Needed for the Proofs

Recall from the notation section that for a sequence of random variables \mathbf{Z}_t , we use the notation $\mathbb{E}_{t-1}[\mathbf{Z}_t]$ to mean the expectation of \mathbf{Z}_t conditioned on all of the previous \mathbf{Z}_t 's. That is:

$$\mathbb{E}_{t-1}[\mathbf{Z}_t] := \mathbb{E}[\mathbf{Z}_t | \mathbf{Z}_1, \dots, \mathbf{Z}_{t-1}]$$

and

$$\mathbb{E}_{t-1}[\mathbf{Z}_t | X] := \mathbb{E}[\mathbf{Z}_t | X, \mathbf{Z}_1, \dots, \mathbf{Z}_{t-1}].$$

Lemma 3.6.1. For $j = 1, \dots, J$ and $k = 1, \dots, K$, for all $X_{\hat{u}_j+k-1} \in \Gamma_{j,k-1}^{\hat{u}_j}$

1. $\mathbf{0} \preceq \frac{1}{\alpha} \sum_{t \in \mathcal{J}_{\hat{u}_j+k}} \mathbb{E}_{t-1} [\mathbf{a}_{t,*} \mathbf{a}_{t,*}' \mid X_{\hat{u}_j+k-1}] = b^2 \mathbf{a}_{t-1,*} \mathbf{a}_{t-1,*}' + (\mathbf{\Lambda}_{\nu,t})_* \preceq (b^2 r \gamma^2 + (1 - b^2) \lambda^+) \mathbf{I}$
2. $(1-b^2) \lambda_{\text{new}}^- \mathbf{I} \preceq \frac{1}{\alpha} \sum_{t \in \mathcal{J}_{\hat{u}_j+k}} \mathbb{E}_{t-1} [\mathbf{a}_{t,\text{new}} \mathbf{a}_{t,\text{new}}' \mid X_{\hat{u}_j+k-1}] = b^2 \mathbf{a}_{t-1,\text{new}} \mathbf{a}_{t-1,\text{new}}' + (\mathbf{\Lambda}_{\nu,t})_{\text{new}} \preceq (b^2 c_{\text{new}} \gamma_{\text{new}}^2 + (1 - b^2) \lambda_{\text{new}}^+) \mathbf{I}$
3. $\frac{1}{\alpha} \sum_{t \in \mathcal{J}_{\hat{u}_j+k}} \mathbb{E}_{t-1} [\mathbf{a}_{t,*} \mathbf{a}_{t,\text{new}}' \mid X_{\hat{u}_j+k-1}] = b^2 \mathbf{a}_{t-1,*} \mathbf{a}_{t-1,\text{new}}'$

with $\hat{u}_j = u_j$ or $\hat{u}_j = u_j + 1$. The only reason we need the assumption $X_{\hat{u}_j+k-1} \in \Gamma_{j,k-1}^{\hat{u}_j}$ is to apply Fact 3.5.13 which allows us to use the “ $\mathbf{a}_{t,\text{new}}$ small” bounds.

The same bounds also hold for summation over $t \in \mathcal{J}_{u_j+1}$ when we condition on $X_{u_j} \in \Gamma_{j-1,\text{end}}$.

The proof is given in Section 3.7. We implicitly use Fact 3.5.13 in the proof.

Using $\|\mathbf{a}_{t,*} \mathbf{a}_{t,*}'\|_2 \leq r \gamma^2$, $\|\mathbf{a}_{t,\text{new}} \mathbf{a}_{t,\text{new}}'\|_2 \leq c \gamma_{\text{new}}^2$ (this holds because of Fact 3.5.13) and applying the matrix Azuma inequality (Lemma 3.A.9 in the appendix) to the first two claims above gives the following lemma.

Lemma 3.6.2. For $j = 1, \dots, J$ and $k = 1, \dots, K$,

1. $\mathbb{P} \left(\lambda_{\max} \left(\frac{1}{\alpha} \sum_{t \in \mathcal{J}_{\hat{u}_j+k}} \mathbf{a}_{t,*} \mathbf{a}_{t,*}' \right) \leq b^2 r \gamma^2 + (1 - b^2) \lambda^+ + \epsilon \mid X_{\hat{u}_j+k-1} \right) \geq 1 - r F(\alpha, \epsilon, r \gamma^2)$
2. $\mathbb{P} \left(\lambda_{\min} \left(\frac{1}{\alpha} \sum_{t \in \mathcal{J}_{\hat{u}_j+k}} \mathbf{a}_{t,\text{new}} \mathbf{a}_{t,\text{new}}' \right) \geq (1 - b^2) \lambda_{\text{new}}^- - \epsilon \mid X_{\hat{u}_j+k-1} \right) \geq 1 - r F(\alpha, \epsilon, c_{\text{new}} \gamma_{\text{new}}^2)$
3. $\mathbb{P} \left(\lambda_{\max} \left(\frac{1}{\alpha} \sum_{t \in \mathcal{J}_{\hat{u}_j+k}} \mathbf{a}_{t,\text{new}} \mathbf{a}_{t,\text{new}}' \right) \leq b^2 c_{\text{new}} \gamma_{\text{new}}^2 + (1 - b^2) \lambda_{\text{new}}^+ + \epsilon \mid X_{\hat{u}_j+k-1} \right) \geq 1 - r F(\alpha, \epsilon, c_{\text{new}} \gamma_{\text{new}}^2)$

for all $X_{\hat{u}_j+k-1} \in \Gamma_{j,0}^{\hat{u}_j}$ with $\hat{u}_j = u_j$ or $\hat{u}_j = u_j + 1$.

The same bounds also hold for summation over $t \in \mathcal{J}_{u_j+1}$ when we condition on $X_{u_j} \in \Gamma_{j-1,\text{end}}$.

Lemma 3.6.3.

$$\mathbb{P} \left(\left\| \frac{1}{\alpha} \sum_{t \in \mathcal{J}_{\hat{u}_j+k}} \mathbf{a}_{t,\text{new}} \mathbf{a}_{t,*}' \right\|_2 \leq \frac{b^2}{\alpha(1-b^2)} \sqrt{c_{\text{new}} r \gamma_{\text{new}}} \gamma + 4\epsilon \mid X_{\hat{u}_j+k-1} \right) \geq 1 - 3(r + c_{\text{new}})F(\alpha, \epsilon, 2\sqrt{c_{\text{new}} r \gamma_{\text{new}}}) - (r + c_{\text{new}})F(\alpha, \epsilon, 4\sqrt{c_{\text{new}} r \gamma_{\text{new}}}).$$

for all $X_{\hat{u}_j+k-1} \in \Gamma_{j,0}^{\hat{u}_j}$ with $\hat{u}_j = u_j$ or $\hat{u}_j = u_j + 1$.

The same bounds also hold for summation over $t \in \mathcal{J}_{u_j+1}$ when we condition on $X_{u_j} \in \Gamma_{j-1,\text{end}}$.

The proof is also in Section 3.7. We use Fact 3.5.13 in the proof.

Remark 3.6.4. Whenever Lemma 3.6.2 or 3.6.3 is applied, we set $\epsilon = 0.01c_{\text{new}}\zeta\lambda^-$.

Lemma 3.6.2 follows directly from Lemma 3.6.1 and the matrix Azuma inequality. Lemma 3.6.3 needs more work to get the $\frac{1}{\alpha}$ factor. This is needed because γ can be large, and the cross term $\mathbf{a}_{t,\text{new}} \mathbf{a}_{t,*}'$ will appear later without an appropriately small factor multiplying it.

Remark 3.6.5. It is possible to also bound $\lambda_{\max}(\frac{1}{\alpha} \sum_t \mathbf{a}_{t,\text{new}} \mathbf{a}_{t,\text{new}}')$ and $\lambda_{\max}(\frac{1}{\alpha} \sum_t \mathbf{a}_{t,*} \mathbf{a}_{t,*}')$ using the same approach used in Lemma 3.6.3. This would give a $\frac{1}{\alpha}$ multiplying the $b^2 c_{\text{new}} \gamma_{\text{new}}^2$ terms.

Lemma 3.6.6. Assume that the assumptions of Theorem 3.2.15 hold. Conditioned on $X_{\hat{u}_j+k-1} \in \Gamma_{j,k-1}^{\hat{u}_j}$ for $\hat{u}_j = u_j$ or $\hat{u}_j = u_j + 1$,

$$\|\mathbf{I}_{\mathcal{T}'} \mathbf{D}_{j,\text{new}}\|_2 \leq \kappa_s^+ := .0215 \quad (3.15)$$

for all \mathcal{T} such that $|\mathcal{T}| \leq s$.

Proof. Recall that $\mathbf{D}_{j,\text{new}} = (\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}') \mathbf{P}_{(j),\text{new}}$. Then $\|\mathbf{I}_{\mathcal{T}'} \mathbf{D}_{j,\text{new}}\|_2 = \|\mathbf{I}_{\mathcal{T}'} (\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}') \mathbf{P}_{(j),\text{new}}\|_2 \leq \|\mathbf{I}_{\mathcal{T}'} \mathbf{P}_{(j),\text{new}}\|_2 + \|\hat{\mathbf{P}}_{(j),*}' \mathbf{P}_{(j),\text{new}}\|_2 \leq \kappa_s(\mathbf{P}_{(j),\text{new}}) + \|\hat{\mathbf{P}}_{(j),*}' (\mathbf{I} - \mathbf{P}_{(j),*} \mathbf{P}_{(j),*}') \mathbf{P}_{(j),\text{new}}\|_2 \leq \kappa_s(\mathbf{P}_{(j),\text{new}}) + \zeta_{j,*}$. The event $X_{\hat{u}_j+k-1} \in \Gamma_{j,k-1}^{\hat{u}_j}$ implies that $\zeta_{j,*} \leq \zeta_{j,*}^+ \leq 0.0015$. Thus, the lemma follows. \square

3.6.2 Proofs of Lemmas 3.5.20, 3.5.21, 3.5.22

Definition 3.6.7. Define the following

1. $\hat{\mathbf{P}}_{(j),\text{new},0} = [\cdot]$ (empty matrix)
2. $\mathbf{D}_{j,\text{new},k} := (\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}' - \hat{\mathbf{P}}_{(j),\text{new},k} \hat{\mathbf{P}}_{(j),\text{new},k}') \mathbf{P}_{(j),\text{new}}$. Recall that $\mathbf{D}_{j,\text{new}} := \mathbf{D}_{j,\text{new},0}$.
3. $\mathbf{D}_{j,*,k} := (\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}' - \hat{\mathbf{P}}_{(j),\text{new},k} \hat{\mathbf{P}}_{(j),\text{new},k}') \mathbf{P}_{(j),*}$ and $\mathbf{D}_{j,*} := \mathbf{D}_{j,*,0}$
4. Recall that $\zeta_{j,0} = \|\mathbf{D}_{j,\text{new}}\|_2$, $\zeta_{j,\text{new},k} = \|\mathbf{D}_{j,\text{new},k}\|_2$, $\zeta_{j,*} = \|\mathbf{D}_{j,*}\|_2$. Also, clearly, $\|\mathbf{D}_{j,*,k}\|_2 \leq \zeta_{j,*}$.

Remark 3.6.8. In the rest of this section, for ease of notation, we do the following.

- We remove the subscript j from $\mathbf{D}_{j,\text{new},k}$ etc., $\mathbf{E}_{j,\text{new}}$ etc. and $\zeta_{j,\text{new},k}$ etc. (from everything in Definitions 3.5.3 and 3.5.16).
- Similarly we also let $X_k := X_{\hat{u}_j+k}$ and $\Gamma_k := \Gamma_{j,k}^{\hat{u}_j}$. Thus, if we say $\mathbb{P}(\text{event} | X_{k-1} \in \Gamma_{k-1}) \geq p_0$ we mean that $\mathbb{P}(\text{event} | X_{u_j+k-1} \in \Gamma_{j,k-1}^{u_j}) \geq p_0$ and $\mathbb{P}(\text{event} | X_{u_j+1+k-1} \in \Gamma_{j,k-1}^{u_j+1}) \geq p_0$.
- Finally, \sum_t refers to $\sum_{t \in \mathcal{J}_u}$ for $u = \hat{u}_j + k$, $k = 1, 2, \dots, K$, $j = 1, 2, \dots, J$.

Also, note the following.

- Recall that \mathcal{T}_t is included in the definition of X_{k-1} , so conditioned on X_{k-1} , the \mathcal{T}_t are deterministic.
- The proof for the bound on \mathbf{A}_u for $u = u_j + 1$ is the same as that for $u = \hat{u}_j + 1$ since in both cases $\hat{\mathbf{P}}_{t,*} = \hat{\mathbf{P}}_{(j),*} = \hat{\mathbf{P}}_{(j-1)}$ and $\hat{\mathbf{P}}_{t,\text{new}} = [\cdot]$ for all $t \in \mathcal{J}_u$. The same is true for the bounds on $\mathbf{A}_{u_j+1,\perp}$ and \mathcal{H}_{u_j+1} .

Lemma 3.6.9. When $X_{k-1} \in \Gamma_{k-1}$,

1. $\|\mathbf{D}_{*,k-1}\|_2 \leq \zeta_{j,*}^+$ for $k = 1, \dots, K$.
2. $\|\mathbf{D}_{\text{new},k-1}\|_2 \leq \zeta_{k-1,\text{new}}^+$ for $k = 1, \dots, K$ (by definition of Γ_{k-1})

$$3. \quad \|[(\Phi_t)_{\mathcal{T}_t}'(\Phi_t)_{\mathcal{T}_t}]^{-1}\|_2 \leq \phi^+ \text{ (from Lemma 3.5.10)}$$

$$4. \quad \lambda_{\min}(\mathbf{R}_{\text{new}}\mathbf{R}_{\text{new}}') \geq 1 - (\zeta_{j,*}^+)^2$$

Proof of Lemma 3.5.20. We obtain the bounds on \mathbf{A}_u for $u = \hat{u}_j + k$ for $k = 1, 2, \dots, K$ and $\hat{u}_j = u_j$ or $u_j + 1$. For $u = \hat{u}_j + k$, recall that $\mathbf{A}_u := \frac{1}{\alpha} \sum_{t \in \mathcal{J}_u} \mathbf{E}_{j,\text{new}}'(\mathbf{I} - \hat{\mathbf{P}}_{(j),*}\hat{\mathbf{P}}_{(j),*}')\ell_t\ell_t'(\mathbf{I} - \hat{\mathbf{P}}_{(j),*}\hat{\mathbf{P}}_{(j),*}')\mathbf{E}_{j,\text{new}}$.

Notice that $\mathbf{E}_{\text{new}}'(\mathbf{I} - \hat{\mathbf{P}}_{(j),*}\hat{\mathbf{P}}_{(j),*}')\ell_t = \mathbf{R}_{\text{new}}\mathbf{a}_{t,\text{new}} + \mathbf{E}_{\text{new}}'\mathbf{D}_*\mathbf{a}_{t,*}$. Let $\mathbf{Z}_t = \mathbf{R}_{\text{new}}\mathbf{a}_{t,\text{new}}\mathbf{a}_{t,\text{new}}'\mathbf{R}_{\text{new}}'$, and let $\mathbf{Y}_t = \mathbf{R}_{\text{new}}\mathbf{a}_{t,\text{new}}\mathbf{a}_{t,*}'\mathbf{D}_*\mathbf{E}_{\text{new}} + \mathbf{E}_{\text{new}}'\mathbf{D}_*\mathbf{a}_{t,*}\mathbf{a}_{t,\text{new}}'\mathbf{R}_{\text{new}}'$, then

$$\mathbf{A}_u \succeq \frac{1}{\alpha} \sum_t \mathbf{Z}_t + \frac{1}{\alpha} \sum_t \mathbf{Y}_t \quad (3.16)$$

Consider $\sum_t \mathbf{Z}_t = \sum_t \mathbf{R}_{\text{new}}\mathbf{a}_{t,\text{new}}\mathbf{a}_{t,\text{new}}'\mathbf{R}_{\text{new}}'$. With probability 1, $\|\mathbf{Z}_t\|_2 \leq c_{\text{new}}\gamma_{\text{new}}^2$. Using a theorem of Ostrowski [27, Theorem 4.5.9], $\lambda_{\min}(\mathbf{Z}_t) = \lambda_{\min}(\mathbf{R}_{\text{new}}\mathbf{a}_{t,\text{new}}\mathbf{a}_{t,\text{new}}'\mathbf{R}_{\text{new}}') \geq \lambda_{\min}(\mathbf{R}_{\text{new}}\mathbf{R}_{\text{new}}')\lambda_{\min}(\mathbf{a}_{t,\text{new}}\mathbf{a}_{t,\text{new}}')$.

Conditioned on X_{k-1} , the matrix \mathbf{R}_{new} is a constant. Using Lemma 3.6.2,

$$\mathbb{P}\left(\lambda_{\min}\left(\frac{1}{\alpha} \sum_t \mathbf{Z}_t\right) \geq (1 - (\zeta_{j,*}^+)^2) [(1 - b^2)\lambda_{\text{new}}^- - \epsilon] \mid X_{k-1}\right) \geq 1 - c_{\text{new}}F(\alpha, \epsilon, c_{\text{new}}\gamma_{\text{new}}^2). \quad (3.17)$$

for all $X_{k-1} \in \Gamma_{k-1}$.

Consider $\mathbf{Y}_t = \mathbf{R}_{\text{new}}\mathbf{a}_{t,\text{new}}\mathbf{a}_{t,*}'\mathbf{D}_*\mathbf{E}_{\text{new}} + \mathbf{E}_{\text{new}}'\mathbf{D}_*\mathbf{a}_{t,*}\mathbf{a}_{t,\text{new}}'\mathbf{R}_{\text{new}}'$. By Lemma 3.6.3

$$\begin{aligned} \mathbb{P}\left(\lambda_{\min}\left(\frac{1}{\alpha} \sum_t \mathbf{Y}_t\right) \geq -2\zeta_{j,*}^+ \left(\frac{b^2}{\alpha(1-b^2)}\sqrt{c_{\text{new}}r}\gamma_{\text{new}}\gamma + 4\epsilon\right) \mid X_{k-1}\right) \geq \\ 1 - 3(r + c_{\text{new}})F(\alpha, \epsilon, 2\sqrt{c_{\text{new}}r}\gamma_{\text{new}}) - (r + c_{\text{new}})F(\alpha, \epsilon, 4\sqrt{c_{\text{new}}r}\gamma_{\text{new}}) \end{aligned} \quad (3.18)$$

for all $X_{k-1} \in \Gamma_{k-1}$. Combining (3.17) and (3.18) the lemma follows. \square

Proof of Lemma 3.5.21. Remark 3.6.8 applies.

We obtain the bounds on $\mathbf{A}_{u,\perp}$ for $u = \hat{u}_j + k$ for $k = 1, 2, \dots, K$ and $\hat{u}_j = u_j$ or $u_j + 1$. For $u = \hat{u}_j + k$, recall that $\mathbf{A}_{u,\perp} := \frac{1}{\alpha} \sum_t \mathbf{E}_{\text{new},\perp}'(\mathbf{I} - \hat{\mathbf{P}}_{(j),*}\hat{\mathbf{P}}_{(j),*}')\ell_t\ell_t'(\mathbf{I} - \hat{\mathbf{P}}_{(j),*}\hat{\mathbf{P}}_{(j),*}')\mathbf{E}_{\text{new},\perp}$. By their definitions, $\mathbf{E}_{\text{new},\perp}'(\mathbf{I} - \hat{\mathbf{P}}_{(j),*}\hat{\mathbf{P}}_{(j),*}')\ell_t = \mathbf{E}_{\text{new},\perp}'\mathbf{D}_*\mathbf{a}_{t,*}$. Thus, $\mathbf{A}_{u,\perp} = \frac{1}{\alpha} \sum_t \mathbf{Z}_t$ with

$\mathbf{Z}_t = \mathbf{E}_{\text{new},\perp}' \mathbf{D}_* \mathbf{a}_{t,*} \mathbf{a}_{t,*}' \mathbf{D}_* \mathbf{E}_{\text{new},\perp}$. Conditioned on $X_{k-1} \in \Gamma_{k-1}$, $\|\mathbf{Z}_t\|_2 \leq (\zeta_{j,*}^+)^2 r \gamma^2$. Using Lemma 3.6.2, we get that

$$\mathbb{P} \left(\lambda_{\max} \left(\frac{1}{\alpha} \sum_t \mathbf{A}_{u,\perp} \right) \leq (\zeta_{j,*}^+)^2 (b^2 r \gamma^2 + (1-b^2) \lambda^+ + \epsilon) \mid X_{k-1} \right) \geq 1 - r F \left(\alpha, \epsilon, (\zeta_{j,*}^+)^2 r \gamma^2 \right) \quad (3.19)$$

for all $X_{k-1} \in \Gamma_{k-1}$. \square

Proof of Lemma 3.5.22. Remark 3.6.8 applies.

Consider the \mathcal{H}_u term. For ease of notation, define

$$\tilde{\ell}_t = (\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}') \ell_t.$$

Using the expression for \mathcal{H}_u given in Definition 3.5.16, and noting that for a basis matrix \mathbf{E} , $\mathbf{E}\mathbf{E}' + \mathbf{E}_\perp \mathbf{E}_\perp' = \mathbf{I}$ we get that

$$\begin{aligned} \mathcal{H}_u = & \frac{1}{\alpha} \sum_{t \in \mathcal{J}_u} \left((\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}') \mathbf{e}_t \mathbf{e}_t' (\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}') \right. \\ & \left. - (\tilde{\ell}_t \mathbf{e}_t' (\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}') + (\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}') \mathbf{e}_t \tilde{\ell}_t') + (\mathbf{F}_t + \mathbf{F}_t') \right) \end{aligned}$$

where

$$\mathbf{F}_t = \mathbf{E}_{\text{new},\perp} \mathbf{E}_{\text{new},\perp}' \tilde{\ell}_t \tilde{\ell}_t' \mathbf{E}_{\text{new}} \mathbf{E}_{\text{new}}'.$$

Thus,

$$\|\mathcal{H}_u\|_2 \leq \left\| \frac{1}{\alpha} \sum_t \mathbf{e}_t \mathbf{e}_t' \right\|_2 + 2 \left\| \frac{1}{\alpha} \sum_t \tilde{\ell}_t \mathbf{e}_t' \right\|_2 + 2 \left\| \frac{1}{\alpha} \sum_t \mathbf{F}_t \right\|_2 \quad (3.20)$$

Next we obtain high probability bounds on each of the three terms on the right hand side of (3.20).

Consider $\|\frac{1}{\alpha} \sum_t \mathbf{e}_t \mathbf{e}_t'\|_2$. Using Lemma 3.5.10, \mathbf{e}_t satisfies (3.9) with probability one for all $X_{k-1} \in \Gamma_{k-1}$. Expanding this expression gives

$$\mathbf{e}_t \mathbf{e}_t' = \left(\mathbf{I}_{\mathcal{T}_t} [(\Phi_t)_{\mathcal{T}_t}' (\Phi_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \Phi_t \ell_t \right) \left(\mathbf{I}_{\mathcal{T}_t} [(\Phi_t)_{\mathcal{T}_t}' (\Phi_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \Phi_t \ell_t \right)'$$

From Lemma 3.5.10, conditioned on $X_{k-1} \in \Gamma_{k-1}$, $\|\mathbf{e}_t \mathbf{e}_t'\|_2 \leq \left[\phi^+ \left(\zeta_{j,*}^+ \sqrt{r} \gamma + \sqrt{c_{\text{new}}} \gamma_{\text{new}} \right) \right]^2$.

This is a looser but simpler bound obtaining by using $\zeta_{j,\text{new},k-1}^+ \leq 1$.

We can further decompose $\mathbf{e}_t \mathbf{e}_t'$ as

$$\mathbf{e}_t \mathbf{e}_t' = \mathbf{term}_{1,t} + \mathbf{term}_{2,t} + \mathbf{term}_{3,t} + \mathbf{term}_{3,t}'$$

where

$$\begin{aligned} \mathbf{term}_{1,t} &= \mathbf{I}_{\mathcal{T}_t} [(\Phi_t)_{\mathcal{T}_t}' (\Phi_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \mathbf{D}_{*,k-1} \mathbf{a}_{t,*} \mathbf{a}_{t,*}' \mathbf{D}_{*,k-1}' \mathbf{I}_{\mathcal{T}_t} [(\Phi_t)_{\mathcal{T}_t}' (\Phi_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \\ \mathbf{term}_{2,t} &= \mathbf{I}_{\mathcal{T}_t} [(\Phi_t)_{\mathcal{T}_t}' (\Phi_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \mathbf{D}_{\text{new},k-1} \mathbf{a}_{t,\text{new}} \mathbf{a}_{t,\text{new}}' \mathbf{D}_{\text{new},k-1}' \mathbf{I}_{\mathcal{T}_t} [(\Phi_t)_{\mathcal{T}_t}' (\Phi_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \\ \mathbf{term}_{3,t} &= \mathbf{I}_{\mathcal{T}_t} [(\Phi_t)_{\mathcal{T}_t}' (\Phi_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \mathbf{D}_{*,k-1} \mathbf{a}_{t,*} \mathbf{a}_{t,\text{new}}' \mathbf{D}_{\text{new},k-1}' \mathbf{I}_{\mathcal{T}_t} [(\Phi_t)_{\mathcal{T}_t}' (\Phi_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \end{aligned}$$

Then by Lemma 3.6.1,

$$\begin{aligned} \mathbb{E}_{t-1}[\mathbf{term}_{1,t} | X_{k-1}] &= \mathbf{I}_{\mathcal{T}_t} [(\Phi_t)_{\mathcal{T}_t}' (\Phi_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \\ &\quad \mathbf{D}_{*,k-1} (b^2 \mathbf{a}_{t-1,*} \mathbf{a}_{t-1,*}' + (\Lambda_{\nu,t})_*) \mathbf{D}_{*,k-1}' \mathbf{I}_{\mathcal{T}_t} [(\Phi_t)_{\mathcal{T}_t}' (\Phi_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \\ \mathbb{E}_{t-1}[\mathbf{term}_{2,t} | X_{k-1}] &= \mathbf{I}_{\mathcal{T}_t} [(\Phi_t)_{\mathcal{T}_t}' (\Phi_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \mathbf{D}_{\text{new},k-1} (b^2 \mathbf{a}_{t-1,\text{new}} \mathbf{a}_{t-1,\text{new}}' + (\Lambda_{\nu,t})_{\text{new}}) \\ &\quad \mathbf{D}_{\text{new},k-1}' \mathbf{I}_{\mathcal{T}_t} [(\Phi_t)_{\mathcal{T}_t}' (\Phi_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \\ \mathbb{E}_{t-1}[\mathbf{term}_{3,t} | X_{k-1}] &= \mathbf{I}_{\mathcal{T}_t} [(\Phi_t)_{\mathcal{T}_t}' (\Phi_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \\ &\quad \mathbf{D}_{*,k-1} b^2 \mathbf{a}_{t-1,*} \mathbf{a}_{t-1,\text{new}}' \mathbf{D}_{\text{new},k-1}' \mathbf{I}_{\mathcal{T}_t} [(\Phi_t)_{\mathcal{T}_t}' (\Phi_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \end{aligned}$$

The following uses Signal Model 3.3.1, Lemma 3.3.3, and Lemma 3.6.1. When $k = 1$ we use Lemma 3.6.6 which gives the bound $\|\mathbf{D}_{\text{new},0}' \mathbf{I}_{\mathcal{T}_t}\|_2 \leq \kappa_s^+$:

$$\begin{aligned} \left\| \frac{1}{\alpha} \sum_t \mathbb{E}_{t-1}[\mathbf{term}_{1,t} | X_{k-1}] \right\|_2 &\leq \rho^2 h^+ (\phi^+)^2 (\zeta_{j,*}^+)^2 (b^2 r \gamma^2 + (1 - b^2) \lambda^+) \\ \left\| \frac{1}{\alpha} \sum_t \mathbb{E}_{t-1}[\mathbf{term}_{2,t} | X_{k-1}] \right\|_2 &\leq \rho^2 h^+ (\phi^+)^2 (\kappa_s^+)^2 (b^2 c_{\text{new}} \gamma_{\text{new}}^2 + (1 - b^2) \lambda_{\text{new}}^+) \\ \left\| \frac{1}{\alpha} \sum_t \mathbb{E}_{t-1}[\mathbf{term}_{3,t} | X_{k-1}] \right\|_2 &\leq \rho^2 h^+ (\phi^+)^2 \kappa_s^+ \zeta_{j,*}^+ b^2 \sqrt{r c_{\text{new}}} \gamma \gamma_{\text{new}}. \end{aligned}$$

And when $k \geq 2$,

$$\begin{aligned} \left\| \frac{1}{\alpha} \sum_t \mathbb{E}_{t-1}[\mathbf{term}_{1,t} | X_{k-1}] \right\|_2 &\leq \rho^2 h^+ (\phi^+)^2 (\zeta_{j,*}^+)^2 (b^2 r \gamma^2 + (1 - b^2) \lambda^+) \\ \left\| \frac{1}{\alpha} \sum_t \mathbb{E}_{t-1}[\mathbf{term}_{2,t} | X_{k-1}] \right\|_2 &\leq \rho^2 h^+ (\phi^+)^2 (\zeta_{j,\text{new},k-1}^+)^2 (b^2 c_{\text{new}} \gamma_{\text{new}}^2 + (1 - b^2) \lambda_{\text{new}}^+) \\ \left\| \frac{1}{\alpha} \sum_t \mathbb{E}_{t-1}[\mathbf{term}_{3,t} | X_{k-1}] \right\|_2 &\leq \rho^2 h^+ (\phi^+)^2 \zeta_{j,*}^+ \zeta_{j,\text{new},k-1}^+ b^2 \sqrt{r c_{\text{new}}} \gamma \gamma_{\text{new}}. \end{aligned}$$

Thus by Lemma 3.A.10 (Azuma Corollary)

$$\mathbb{P}\left(\left\|\frac{1}{\alpha}\sum_t \mathbf{e}_t \mathbf{e}_t'\right\|_2 \leq b_{2,k} + \epsilon \mid X_{k-1}\right) \geq 1 - nF\left(\alpha, \epsilon, 2\left[\phi^+\left(\zeta_{j,*}^+ \sqrt{r}\gamma + \sqrt{c_{\text{new}}}\gamma_{\text{new}}\right)\right]^2\right) \quad (3.21)$$

for all $X_{k-1} \in \Gamma_{k-1}$.

Next, consider $\left\|\frac{1}{\alpha}\sum_t \tilde{\ell}_t \mathbf{e}_t'\right\|_2$. Observe that when $X_{k-1} \in \Gamma_{k-1}$,

$$\begin{aligned} \tilde{\ell}_t \mathbf{e}_t' &= (\mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}') (\mathbf{P}_* \mathbf{a}_{t,*} + \mathbf{P}_{\text{new}} \mathbf{a}_{t,\text{new}}) (\mathbf{P}_* \mathbf{a}_{t,*} + \mathbf{P}_{\text{new}} \mathbf{a}_{t,\text{new}})' \Phi_t' \mathbf{I}_{\mathcal{T}_t} [(\Phi_t)_{\mathcal{T}_t}' (\Phi_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \\ &= (\mathbf{D}_* \mathbf{a}_{t,*} + \mathbf{D}_{\text{new}} \mathbf{a}_{t,\text{new}}) (\mathbf{D}_{*,k-1} \mathbf{a}_{t,*} + \mathbf{D}_{\text{new},k-1} \mathbf{a}_{t,\text{new}})' \mathbf{I}_{\mathcal{T}_t} [(\Phi_t)_{\mathcal{T}_t}' (\Phi_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \end{aligned}$$

Redefine

$$\begin{aligned} \mathbf{term}_{1,t} &:= (\mathbf{D}_* \mathbf{a}_{t,*} \mathbf{a}_{t,*}' \mathbf{D}_{*,k-1}' + \mathbf{D}_{\text{new}} \mathbf{a}_{t,\text{new}} \mathbf{a}_{t,\text{new}}' \mathbf{D}_{\text{new},k-1}') \mathbf{I}_{\mathcal{T}_t} [(\Phi_t)_{\mathcal{T}_t}' (\Phi_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \\ \mathbf{term}_{2,t} &:= (\mathbf{D}_* \mathbf{a}_{t,*} \mathbf{a}_{t,\text{new}}' \mathbf{D}_{\text{new},k-1}' + \mathbf{D}_{\text{new}} \mathbf{a}_{t,\text{new}} \mathbf{a}_{t,*}' \mathbf{D}_{*,k-1}') \mathbf{I}_{\mathcal{T}_t} [(\Phi_t)_{\mathcal{T}_t}' (\Phi_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \end{aligned}$$

When $k = 1$ we use Lemma 3.6.6 which gives the bound $\|\mathbf{D}_{\text{new}}' \mathbf{I}_{\mathcal{T}_t}\|_2 \leq \kappa_s^+$. By Lemma 3.6.2,

$$\begin{aligned} \mathbb{P}\left(\left\|\frac{1}{\alpha}\sum_{t \in \mathcal{J}_{\hat{u}_j+1}} \mathbf{term}_{1,t}\right\|_2 \leq \left[(\zeta_{j,*}^+)^2 (b^2 r \gamma^2 + (1-b^2) \lambda^+ + \epsilon) \right. \right. \\ \left. \left. + \zeta_{j,\text{new},k-1}^+ \kappa_s^+ (b^2 c_{\text{new}} \gamma_{\text{new}}^2 + (1-b^2) \lambda_{\text{new}}^+ + \epsilon)\right] \phi^+ \mid X_0\right) \\ \geq 1 - rF(\alpha, \epsilon, r\gamma^2) - cF(\alpha, \epsilon, c_{\text{new}} \gamma_{\text{new}}^2). \end{aligned}$$

for all $X_0 \in \Gamma_0$. And by Lemma 3.6.3,

$$\begin{aligned} \mathbb{P}\left(\left\|\frac{1}{\alpha}\sum_{t \in \mathcal{J}_{\hat{u}_j+1}} \mathbf{term}_{2,t}\right\|_2 \leq \left[\zeta_{j,*}^+ \zeta_{j,\text{new},k-1}^+ \kappa_s^+ \left(\frac{b^2}{\alpha(1-b^2)} \sqrt{c_{\text{new}}} r \gamma_{\text{new}} \gamma + 4\epsilon\right) \right. \right. \\ \left. \left. + \zeta_{j,*}^+ \left(\frac{b^2}{\alpha(1-b^2)} \sqrt{c_{\text{new}}} r \gamma_{\text{new}} \gamma + 4\epsilon\right)\right] \phi^+ \mid X_0\right) \\ \geq 1 - 3(r + c_{\text{new}})F(\alpha, \epsilon, 2\sqrt{c_{\text{new}}} r \gamma_{\text{new}}) - (r + c_{\text{new}})F(\alpha, \epsilon, 4\sqrt{c_{\text{new}}} r \gamma_{\text{new}}) \end{aligned}$$

for all $X_0 \in \Gamma_0$.

When $k \geq 2$ we apply Cauchy-Schwarz (Lemma 3.A.3) to $\mathbf{term}_{1,t}$. Here $\mathbf{X}_t = (\mathbf{D}_* \mathbf{a}_{t,*} \mathbf{a}_{t,*}' \mathbf{D}_{*,k-1}' + \mathbf{D}_{\text{new}} \mathbf{a}_{t,\text{new}} \mathbf{a}_{t,\text{new}}' \mathbf{D}_{\text{new},k-1}')$ and $\mathbf{Y}_t = \mathbf{I}_{\mathcal{T}_t} [(\Phi_t)_{\mathcal{T}_t}' (\Phi_t)_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}'$.

We can bound the norm of $\frac{1}{\alpha} \sum_t \mathbf{X}_t \mathbf{X}_t'$ using Lemmas 3.6.2 and 3.6.9. By Lemma 3.3.3 we get

$$\left\| \frac{1}{\alpha} \sum_t \mathbf{Y}_t \mathbf{Y}_t' \right\|_2 \leq \rho^2 h^+ (\phi^+)^2$$

So when $k \geq 2$,

$$\begin{aligned} \mathbb{P} \left(\left\| \frac{1}{\alpha} \sum_{t \in \mathcal{J}_{\hat{u}_j+k}} \mathbf{term}_{1,t} \right\|_2 \leq \left[(\zeta_{j,*}^+)^2 (b^2 r \gamma^2 + (1-b^2) \lambda^+ + \epsilon) + \right. \\ \left. \zeta_{j,\text{new},k-1}^+ (b^2 c_{\text{new}} \gamma_{\text{new}}^2 + (1-b^2) \lambda_{\text{new}}^+ + \epsilon) \right] \left(\sqrt{\rho^2 h^+} \phi^+ \right) \mid X_{k-1} \right) \\ \geq 1 - r F(\alpha, \epsilon, r \gamma^2) - c_{\text{new}} F(\alpha, \epsilon, c_{\text{new}} \gamma_{\text{new}}^2). \end{aligned}$$

for all $X_{k-1} \in \Gamma_{k-1}$. Similarly, using Lemmas 3.3.3, 3.6.3, and 3.6.9,

$$\begin{aligned} \mathbb{P} \left(\left\| \frac{1}{\alpha} \sum_{t \in \mathcal{J}_{\hat{u}_j+k}} \mathbf{term}_{2,t} \right\|_2 \leq \left[\zeta_{j,*}^+ \zeta_{j,\text{new},k-1}^+ \left(\frac{b^2}{\alpha(1-b^2)} \sqrt{c_{\text{new}} r \gamma_{\text{new}} \gamma} + 4\epsilon \right) + \right. \\ \left. \zeta_{j,*}^+ \left(\frac{b^2}{\alpha(1-b^2)} \sqrt{c_{\text{new}} r \gamma_{\text{new}} \gamma} + 4\epsilon \right) \right] \left(\sqrt{\rho^2 h^+} \phi^+ \right) \mid X_{k-1} \right) \\ \geq 1 - 3(r + c_{\text{new}}) F(\alpha, \epsilon, 2\sqrt{c_{\text{new}} r \gamma_{\text{new}} \gamma}) - (r + c_{\text{new}}) F(\alpha, \epsilon, 4\sqrt{c_{\text{new}} r \gamma_{\text{new}} \gamma}) \end{aligned}$$

for all $X_{k-1} \in \Gamma_{k-1}$.

Thus,

$$\begin{aligned} \mathbb{P} \left(\left\| \frac{1}{\alpha} \sum_t \tilde{\ell}_t \mathbf{e}_t' \right\|_2 \leq b_{4,k} \mid X_{k-1} \right) \geq 1 - r F(\alpha, \epsilon, r \gamma^2) - c_{\text{new}} F(\alpha, \epsilon, c_{\text{new}} \gamma_{\text{new}}^2) - \\ 3(r + c_{\text{new}}) F(\alpha, \epsilon, 2\sqrt{c_{\text{new}} r \gamma_{\text{new}} \gamma}) - (r + c_{\text{new}}) F(\alpha, \epsilon, 4\sqrt{c_{\text{new}} r \gamma_{\text{new}} \gamma}) \end{aligned} \quad (3.22)$$

for all $X_{k-1} \in \Gamma_{k-1}$.

Finally, consider $\left\| \frac{1}{\alpha} \sum_t \mathbf{F}_t \right\|_2$. Notice that

$$\begin{aligned} \mathbf{F}_t &= \mathbf{E}_{\text{new},\perp} \mathbf{E}_{\text{new},\perp}' \tilde{\ell}_t \tilde{\ell}_t' \mathbf{E}_{\text{new}} \mathbf{E}_{\text{new}}' \\ &= \mathbf{E}_{\text{new},\perp} \mathbf{E}_{\text{new},\perp}' (\mathbf{D}_* \mathbf{a}_{t,*} + \mathbf{D}_{\text{new}} \mathbf{a}_{t,\text{new}}) (\mathbf{D}_* \mathbf{a}_{t,*} + \mathbf{D}_{\text{new}} \mathbf{a}_{t,\text{new}})' \mathbf{E}_{\text{new}} \mathbf{E}_{\text{new}}' \\ &= \mathbf{E}_{\text{new},\perp} \mathbf{E}_{\text{new},\perp}' (\mathbf{D}_* \mathbf{a}_{t,*}) (\mathbf{D}_* \mathbf{a}_{t,*} + \mathbf{D}_{\text{new}} \mathbf{a}_{t,\text{new}})' \mathbf{E}_{\text{new}} \mathbf{E}_{\text{new}}' \\ &= \mathbf{E}_{\text{new},\perp} \mathbf{E}_{\text{new},\perp}' (\mathbf{D}_* \mathbf{a}_{t,*} \mathbf{a}_{t,*}' \mathbf{D}_* + \mathbf{D}_* \mathbf{a}_{t,*} \mathbf{a}_{t,\text{new}}' \mathbf{D}_{\text{new}} + \mathbf{D}_{\text{new}} \mathbf{a}_{t,\text{new}} \mathbf{a}_{t,*}' \mathbf{D}_*)' \mathbf{E}_{\text{new}} \mathbf{E}_{\text{new}}' \end{aligned}$$

Then by Lemmas 3.6.2 and 3.6.3,

$$\begin{aligned} \mathbb{P} \left(\left\| \frac{1}{\alpha} \sum_t \mathbf{F}_t \right\|_2 \leq b_6 \mid X_{k-1} \in \Gamma_{k-1} \right) &\geq 1 - c_{\text{new}} F(\alpha, \epsilon, r\gamma^2 | X_{k-1}) - \\ &3(r + c_{\text{new}})F(\alpha, \epsilon, 2\sqrt{c_{\text{new}}r}\gamma\gamma_{\text{new}}) - (r + c_{\text{new}})F(\alpha, \epsilon, 4\sqrt{c_{\text{new}}r}\gamma\gamma_{\text{new}}) \end{aligned} \quad (3.23)$$

for all $X_{k-1} \in \Gamma_{k-1}$.

Using (3.20), (3.21), (3.22) and (3.23) and the union bound, for any $X_{k-1} \in \Gamma_{k-1}$, we get the result. \(\square\)

3.7 Proof of Lemmas 3.6.1 and 3.6.3

Proof of Lemma 3.6.2. For the first claim, begin by observing that $\mathbf{a}_{t,*}\mathbf{a}_{t,*}'$ is positive semidefinite. We also have

$$\begin{aligned} \mathbb{E}_{t-1}[\mathbf{a}_{t,*}\mathbf{a}_{t,*}' | X_{k-1}] &= \mathbb{E}_{t-1}[(b\mathbf{a}_{t-1,*} + \boldsymbol{\nu}_{t,*})(b\mathbf{a}_{t-1,*} + \boldsymbol{\nu}_{t,*})' | X_{k-1}] \\ &= \mathbb{E}_{t-1}[b^2\mathbf{a}_{t-1,*}\mathbf{a}_{t-1,*}' + b\mathbf{a}_{t-1,*}\boldsymbol{\nu}_{t,*}' + b\boldsymbol{\nu}_{t,*}\mathbf{a}_{t-1,*}' + \boldsymbol{\nu}_{t,*}\boldsymbol{\nu}_{t,*}' | X_{k-1}] \\ &= b^2\mathbf{a}_{t-1,*}\mathbf{a}_{t-1,*}' + \boldsymbol{\Lambda}_{\nu,t} \\ &\preceq (b^2r\gamma^2 + (1 - b^2)\lambda^+) \mathbf{I} \end{aligned}$$

The cross terms are zero because of Lemma 3.A.4. Notice that $\boldsymbol{\nu}_{t,*}$ is zero mean and independent of $\boldsymbol{\nu}_{\tau,*}$ for $\tau < t$, and $\mathbf{a}_{1,*}, \dots, \mathbf{a}_{t-1,*}$ are functions of $\{\boldsymbol{\nu}_{\tau,*}\}$ for $\tau = 1, \dots, t-1$.

Claim 2 is done in exact same manner, except that we also need the fact that $\lambda_{\min}(\mathbf{A} + \mathbf{B}) \geq \lambda_{\min}(\mathbf{A}) + \lambda_{\min}(\mathbf{B})$.

Claim 3 uses the same expansion and the fact that $\boldsymbol{\nu}_t$ has diagonal covariance, so $\mathbb{E}[\boldsymbol{\nu}_{t,*}\boldsymbol{\nu}_{t,\text{new}}' | X_{k-1}] = \mathbf{0}$. \(\square\)

Proof of Lemma 3.6.3. By Lemma 3.A.6 with $\mathbf{c}_\tau = \mathbf{a}_{\tau+(\hat{u}_j+k-1)\alpha+1,\text{new}}$, $\boldsymbol{\mu}_\tau = \boldsymbol{\nu}_{\tau+(\hat{u}_j+k-1)\alpha+1,\text{new}}$, $\tilde{\mathbf{c}}_\tau = \mathbf{a}_{\tau+(\hat{u}_j+k-1)\alpha+1,*}$ and $\tilde{\boldsymbol{\mu}}_\tau = \boldsymbol{\nu}_{\tau+(\hat{u}_j+k-1)\alpha+1,*}$ for $\tau = 0, 1, \dots, \alpha-1$, we get

$$\sum_{t \in \mathcal{J}_{\hat{u}_j+k}} \mathbf{a}_{t,\text{new}}\mathbf{a}_{t,*}' = \sum_{i \in \mathcal{J}_{\hat{u}_j+k}} [\mathbf{Z}_{1,i} + \mathbf{Z}_{2,i} + \mathbf{Z}_{3,i} + \mathbf{Z}_{4,i}] + \mathbf{Z}_5$$

with

$$\begin{aligned}
\mathbf{Z}_{1,i} &= \frac{(1 - b^{2(\alpha-i)})}{1 - b^2} \boldsymbol{\nu}_{i,\text{new}} \boldsymbol{\nu}_{i,*}', \\
\mathbf{Z}_{2,i} &= \sum_{i_2=(\hat{u}_j+k-1)\alpha+1}^{i-1} \frac{(1 - b^{2(\alpha-i)})}{1 - b^2} b^{i-i_2} \boldsymbol{\nu}_{i,\text{new}} \boldsymbol{\nu}_{i_2,*}', \\
\mathbf{Z}_{3,i} &= \sum_{i_2=(\hat{u}_j+k)\alpha+1-i}^{(\hat{u}_j+k)\alpha} \frac{(1 - b^{2(\alpha-i_2)})}{1 - b^2} b^{i+i_2-\alpha+1} \boldsymbol{\nu}_{\alpha-i-1,\text{new}} \boldsymbol{\nu}_{i_2,*}', \\
\mathbf{Z}_{4,i} &= \frac{b^{i+1}(1 - b^{2(\alpha-i)})}{1 - b^2} (\boldsymbol{\nu}_{i,\text{new}} \mathbf{a}_{(\hat{u}_j+k-1)\alpha,*}' + \mathbf{a}_{(\hat{u}_j+k-1)\alpha,\text{new}} \boldsymbol{\nu}_{i,*}') \\
\mathbf{Z}_5 &= \frac{b^2(1 - b^{2\alpha})}{1 - b^2} \mathbf{a}_{(\hat{u}_j+k-1)\alpha,\text{new}} \mathbf{a}_{(\hat{u}_j+k-1)\alpha,*}'
\end{aligned}$$

Using $\|\boldsymbol{\nu}_t\|_\infty \leq (1-b)\gamma$, $\|\boldsymbol{\nu}_{t,\text{new}}\|_\infty \leq (1-b)\gamma_{\text{new}}$, and the geometric series formula, we get the following norm bounds (recall that $b < 1$):

1. $\|\mathbf{Z}_{1,i}\|_2 \leq \frac{(1 - b^{2\alpha})(1 - b)}{1 + b} \sqrt{c_{\text{new}} r} \gamma_{\text{new}} \gamma \leq \sqrt{c_{\text{new}} r} \gamma_{\text{new}} \gamma$
2. $\|\mathbf{Z}_{2,i}\|_2 \leq \sqrt{c_{\text{new}} r} \gamma_{\text{new}} \gamma$
3. $\|\mathbf{Z}_{3,i}\|_2 \leq \sqrt{c_{\text{new}} r} \gamma_{\text{new}} \gamma$
4. $\|\mathbf{Z}_{4,i}\|_2 \leq 2\sqrt{c_{\text{new}} r} \gamma_{\text{new}} \gamma$
5. $\|\mathbf{Z}_5\|_2 \leq \frac{b^2}{1 - b^2} \sqrt{c_{\text{new}} r} \gamma_{\text{new}} \gamma$

The bounds 1-4 may be somewhat loose, but they are only needed to lower bound the probability of the good event. Bound 5) will appear in the $\zeta_{k,\text{new}}^+$ expression, so we retain the b 's.

In all expectations, we need to condition on X_{k-1} for all $X_{k-1} \in \Gamma_{k-1}$ in order use Fact 3.5.13. This is necessary to ensure that the tighter bound γ_{new} applies for the given time interval. By diagonal covariance $\mathbb{E}_{i-1}[\mathbf{Z}_{1,i}|X_{k-1}] = \mathbf{0}$. The proof of $\mathbb{E}_{i-1}[\mathbf{Z}_{2,i}] = \mathbb{E}_{i-1}[\mathbf{Z}_{3,i}] = \mathbf{0}$ is shown below. $\mathbb{E}_{i-1}[\mathbf{Z}_{4,i}] = \mathbf{0}$ because the $\boldsymbol{\nu}_t$ are zero-mean and $\mathbf{a}_{(\hat{u}_j+k-1)\alpha}$ is constant conditioned on X_{k-1} . \mathbf{Z}_5 will be the non-zero term.

Ignoring the scalar coefficients,

$$\begin{aligned}
\mathbb{E}_{i-1}[\mathbf{Z}_{2,i}] &= \mathbb{E}_{i-1} \left[\sum_{i_2=(\hat{u}_j+k-1)\alpha+1}^{i-1} \boldsymbol{\nu}_{i,\text{new}} \boldsymbol{\nu}_{i_2,*}' \right] \\
&= \mathbb{E} \left[\boldsymbol{\nu}_{i,\text{new}} \sum_{i_2=(\hat{u}_j+k-1)\alpha+1}^{i-1} \boldsymbol{\nu}_{i_2,*}' \mid \right. \\
&\quad \sum_{i_2=(\hat{u}_j+k-1)\alpha+1}^{i-2} \boldsymbol{\nu}_{i-1,\text{new}} \boldsymbol{\nu}_{i_2,*}', \\
&\quad \left. \sum_{i_2=(\hat{u}_j+k-1)\alpha+1}^{i-3} \boldsymbol{\nu}_{i-2,\text{new}} \boldsymbol{\nu}_{i_2,*}', \dots, \boldsymbol{\nu}_{(\hat{u}_j+k-1)\alpha+2,\text{new}} \boldsymbol{\nu}_{(\hat{u}_j+k-1)\alpha+1,*}' \right] \\
&= \mathbf{0}
\end{aligned}$$

Notice that everything else in the above expression is independent of $\boldsymbol{\nu}_{i,\text{new}}$, so by Lemma 3.A.4 the expectation is zero.

Similarly (again ignore scalar coefficients for simplicity) and letting $\alpha \equiv (\hat{u}_j + k)\alpha + 1$

$$\begin{aligned}
\mathbb{E}_{i-1}[\mathbf{Z}_{3,i}] &= \mathbb{E}_{i-1} \left[\sum_{i_2=\alpha-i}^{\alpha-1} \boldsymbol{\nu}_{\alpha-i-1,\text{new}} \boldsymbol{\nu}_{i_2,*}' \right] \\
&= \mathbb{E} \left[\boldsymbol{\nu}_{\alpha-i-1,\text{new}} \sum_{i_2=\alpha-i}^{\alpha-1} \boldsymbol{\nu}_{i_2,*}' \mid \right. \\
&\quad \sum_{i_2=\alpha-i+1}^{\alpha-1} \boldsymbol{\nu}_{\alpha-i,\text{new}} \boldsymbol{\nu}_{i_2,*}', \\
&\quad \left. \sum_{i_2=\alpha-i+2}^{\alpha-1} \boldsymbol{\nu}_{\alpha-i+1,\text{new}} \boldsymbol{\nu}_{i_2,*}', \dots, \sum_{i_2=\alpha-1}^{\alpha-1} \boldsymbol{\nu}_{\alpha-2,\text{new}} \boldsymbol{\nu}_{i_2,*}' \right].
\end{aligned}$$

As before, all terms are independent of $\boldsymbol{\nu}_{\alpha-i-1}$ so the expectation is zero by Lemma 3.A.4.

Using the Azuma inequality,

1. $\mathbb{P} \left(\left\| \frac{1}{\alpha} \sum_{i=0}^{\alpha-1} \mathbf{Z}_{1,i} \right\|_2 \leq \epsilon \right) \geq 1 - (r + c_{\text{new}}) F(\alpha, \epsilon, 2\sqrt{c_{\text{new}} r \gamma \gamma_{\text{new}}})$
2. $\mathbb{P} \left(\left\| \frac{1}{\alpha} \sum_{i=0}^{\alpha-1} \mathbf{Z}_{2,i} \right\|_2 \leq \epsilon \right) \geq 1 - (r + c_{\text{new}}) F(\alpha, \epsilon, 2\sqrt{c_{\text{new}} r \gamma \gamma_{\text{new}}})$
3. $\mathbb{P} \left(\left\| \frac{1}{\alpha} \sum_{i=0}^{\alpha-1} \mathbf{Z}_{3,i} \right\|_2 \leq \epsilon \right) \geq 1 - (r + c_{\text{new}}) F(\alpha, \epsilon, 2\sqrt{c_{\text{new}} r \gamma \gamma_{\text{new}}})$
4. $\mathbb{P} \left(\left\| \frac{1}{\alpha} \sum_{i=0}^{\alpha-1} \mathbf{Z}_{4,i} \right\|_2 \leq \epsilon \right) \geq 1 - (r + c_{\text{new}}) F(\alpha, \epsilon, 4\sqrt{c_{\text{new}} r \gamma \gamma_{\text{new}}})$

$$5. \mathbb{P} \left(\left\| \frac{1}{\alpha} \mathbf{Z}_5 \right\|_2 \leq \frac{b^2}{\alpha(1-b^2)} \sqrt{c_{\text{new}} r \gamma_{\text{new}} \gamma} \right) = 1$$

Combining the above yields the claim of Lemma 3.6.3.

□

3.8 Alternative Subspace Model and Algorithm

3.8.1 Deletion Model

Recall Signal Model 3.2.8 which assumes

$$\mathbf{P}_t = \begin{cases} [\mathbf{P}_{t-1} \mathbf{R}_t \setminus \mathbf{P}_{t,\text{old}} \mathbf{P}_{t,\text{new}}] & \text{if } t = t_1 \text{ or } t_2 \text{ or } \dots t_J \\ \mathbf{P}_{t-1} & \text{otherwise} \end{cases} \quad (3.24)$$

According to this model, at the change times t_j , some directions in the span of \mathbf{P}_{t-1} may be removed, and new directions may be added.

In this section we will introduce an extension of Algorithm 3 which deletes directions from the estimate of $\text{span}(\mathbf{P}_t)$ by re-estimating the previous subspace before beginning to estimate new directions by the same projection PCA procedure. We prove a result similar to Theorem 3.2.15 for this algorithm.

Definition 3.8.1. *Define*

- $c_{\text{new}} := \max_j c_{j,\text{new}}$
- $c_{\text{dif}} := \max_j \sum_{i=1}^j (c_{j,\text{new}} - c_{j,\text{old}})$
- $r_j := \text{rank}(\mathbf{P}_{(j)})$
- $r_{\text{max}} := \max_j r_j = r_0 + c_{\text{dif}}$
- $r = r_0 + J c_{\text{new}}$

In order to generate a new estimate of the subspace, we need a clustering assumption on the eigenvalues of $\mathbf{\Lambda}_{a,t}$ after the subspace change has stabilized. The reason for the cluster-PCA algorithm and clustering assumption is that the condition number of $\mathbf{\Lambda}_{a,t}$ may be large, and the error in the PCA step, \mathbf{e}_t , is correlated with the true data ℓ_t [12].

Model 3.8.2. Assume Signal Model 3.2.8 and assume:

1. During the interval $[t_{j+1} - d_2, t_{j+1} - 1]$, $\mathbf{\Lambda}_{a,t}$ is constant
2. There exists a partition $\mathcal{G}_{j,1}, \mathcal{G}_{j,2}, \dots, \mathcal{G}_{j,\vartheta}$ of the index set $\{1, 2, \dots, r_j\}$ with $\min_{i \in \mathcal{G}_{j,k}} \lambda_i(\mathbf{\Lambda}_{a,t_{j+1}-1}) > \max_{i \in \mathcal{G}_{j,k+1}} \lambda_i(\mathbf{\Lambda}_{a,t_{j+1}-1})$ for $k = 1, \dots, \vartheta - 1$ that satisfies

$$\tilde{g} \leq \tilde{g}^+ \quad \text{and} \quad \tilde{\chi} \leq \tilde{\chi}^+ \quad \text{and} \quad \vartheta \leq \vartheta^+$$

for a $\tilde{g}^+ > 1$ but not too large and a $\tilde{\chi}^+ < 1$. Where

$$\tilde{g} := \max_{j,k} \frac{\max_{i \in \mathcal{G}_{j,k}} \lambda_i(\mathbf{\Lambda}_{a,t_{j+1}-d_2})}{\min_{i \in \mathcal{G}_{j,k}} \lambda_i(\mathbf{\Lambda}_{a,t_{j+1}-d_2})} \quad \text{and} \quad \tilde{\chi} := \max_{j,k} \frac{\max_{i \in \mathcal{G}_{j,k+1}} \lambda_i(\mathbf{\Lambda}_{a,t_{j+1}-d_2})}{\min_{i \in \mathcal{G}_{j,k}} \lambda_i(\mathbf{\Lambda}_{a,t_{j+1}-d_2})}.$$

Notice that $\tilde{g} \geq 1$ and measures how close the eigenvalues within a cluster are, and $\tilde{\chi} \leq 1$ and measures how far apart the adjacent clusters are.

Definition 3.8.3. Define

1. $\mathbf{G}_{j,k} := (\mathbf{P}_{(j)})_{\mathcal{G}_{j,k}}$ to be the eigenvectors corresponding to the eigenvalues indexed by $\mathcal{G}_{j,k}$, so $\text{range}(\mathbf{P}_{(j)}) = \text{range}([\mathbf{G}_{j,1} \ \mathbf{G}_{j,2} \ \dots \ \mathbf{G}_{j,\vartheta}])$;
2. $\tilde{c}_{j,k} := |\mathcal{G}_{j,k}| = \text{rank}(\mathbf{G}_{j,k})$, so $\sum_{k=1}^{\vartheta} \tilde{c}_{j,k} = r_j$;
3. $\lambda_{j,k}^- := \min_{i \in \mathcal{G}_{j,k}} \min_{t \in [t_{j+1}-d_2, t_{j+1}-1]} \lambda_i(\mathbf{\Lambda}_{a,t})$ and $\lambda_{j,k}^+ := \max_{i \in \mathcal{G}_{j,k}} \max_{t \in [t_{j+1}-d_2, t_{j+1}-1]} \lambda_i(\mathbf{\Lambda}_{a,t})$
4. $\tilde{g}_{j,k} := \frac{\lambda_{j,k}^+}{\lambda_{j,k}^-}$ and $\tilde{\chi}_{j,k} := \frac{\lambda_{j,k+1}^+}{\lambda_{j,k}^-}$ (notice $\tilde{g} = \max_{j,k} \tilde{g}_{j,k}$ and $\tilde{\chi} = \max_{j,k} \tilde{\chi}_{j,k}$)

3.8.2 Performance Guarantee for Algorithm 4

Theorem 3.8.4 (Correctness result for Algorithm 4 under Signal Model 3.8.2). Pick a ζ that satisfies

$$\zeta \leq \min \left\{ \frac{10^{-4}}{(r_{\max} + c_{\text{new}})^2}, \frac{1.5 \times 10^{-4}}{(r_{\max} + c_{\text{new}})^2 f}, \frac{1}{(r_{\max} + c_{\text{new}})^3 \gamma^2}, \frac{0.01 \lambda^-}{b^2 (r_{\max} + c_{\text{new}})^3 \gamma^2}, \gamma \right\}.$$

Suppose

1. $\|(I - \hat{\mathbf{P}}_{(0)} \hat{\mathbf{P}}_{(0)}') \mathbf{P}_{(0)}\|_2 \leq r_0 \zeta$;

2. The algorithm parameters are set as:

- $K = \left\lceil \frac{\log(0.16c_{\text{new}}\zeta)}{\log(0.4)} \right\rceil$;
- $\xi = \sqrt{c_{\text{new}}}\gamma_{\text{new}} + (\sqrt{r} + \sqrt{c_{\text{new}}})\sqrt{\zeta}$;
- $\omega = 7\xi$;
- $\alpha = C_1(\log(6KJ) + 11\log(n))$ for a constant $C_1 \geq C_{\text{add}}$ with

$$C_{\text{add}} := \frac{4800}{(\zeta\lambda^-)^2} \max\{16, (1.2\xi)^4\}$$

- $\tilde{\alpha} = C_2(\log(50\vartheta J) + 11\log(n))$ for a $C_2 \geq C_{\text{del}}$ with

$$C_{\text{del}} := \frac{8 \cdot 100^2 \cdot 4^2 \cdot r^2 \gamma^4}{(c_{\text{new}}\zeta\lambda^-)^2}$$

3. Signal Model 3.8.2 holds with $b \leq 0.1$ and

- $d_1 \geq (K+2)\alpha$ and $d_2 \geq \vartheta^+\tilde{\alpha} + 2\alpha$;
- $t_{j+1} - t_j \geq d_1 + d_2$ for all j ;
- $\sqrt{c_{\text{new}}}\gamma_{\text{new}} + (\sqrt{r} + \sqrt{c_{\text{new}}})\sqrt{\zeta} \leq \frac{x_{\min}}{14}$;
- $g \leq \sqrt{2}$;
- $b^2 c_{\text{new}} \eta_{\text{new}} g \leq 0.5$ (as before, because $b \leq 0.1$, this will be satisfied if $c_{\text{new}} \eta_{\text{new}} g \leq 50$)

4. The clustering assumption holds with $\tilde{g}^+ = 1.5$, $\tilde{\chi}^+ = 0.2$, and $\vartheta^+ = 3$.

5. The support of \mathbf{x}_t changes enough so that for α and $\tilde{\alpha}$ as chosen above, Signal Model 3.2.11 holds with $\beta = h^+ \min\{\alpha, \tilde{\alpha}\}$ and

$$\lceil \varrho \rceil^2 h^+ \leq 0.0024, \text{ and } \varrho_2 s \max\{\alpha, \tilde{\alpha}\} \leq n$$

or Signal Model 3.2.12 holds with $s \leq (6 \times 10^{-4}) \min\{\alpha, \tilde{\alpha}\}$ and $\max\{\alpha, \tilde{\alpha}\} \leq \frac{n}{m}$.

6. The low dimensional subspace is dense such that

- $\max_j \kappa_{2s}(\mathbf{P}_{(j)}) \leq 0.3$;
- $\max_j \kappa_{2s}(\mathbf{P}_{(j),\text{new}}) \leq 0.02$.

Then, with probability at least $1 - 2n^{-10}$, at all times t

1. The support of \mathbf{x}_t is recovered exactly, i.e. $\hat{\mathcal{T}}_t = \mathcal{T}_t$
2. The estimate of the subspace change time satisfies $t_j \leq \hat{t}_j \leq t_j + 2\alpha$, for $j = 1, \dots, J$;
3. The estimate of the number of new directions is correct, i.e. $\hat{c}_{j,\text{new},k} = c_{j,\text{new}}$ for $j = 1, \dots, J$ and $k = 1, \dots, K$;
4. The recovery error satisfies:

$$\|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2 \leq \begin{cases} 1.2(\sqrt{\zeta} + \sqrt{c_{\text{new}}}\gamma_{\text{new}}) & t \in [t_j, \hat{t}_j] \\ 1.2(1.84\sqrt{\zeta} + (0.4)^{k-1}\sqrt{c_{\text{new}}}\gamma_{\text{new}}) & t \in [\hat{t}_j + (k-1)\alpha, \hat{t}_j + k\alpha - 1], \\ & k = 1, 2, \dots, K \\ 2.4\sqrt{\zeta} & t \in [\hat{t}_j + K\alpha, t_{j+1} - 1] \end{cases}$$

5. For $j = 1, \dots, J$, $t_j \leq \hat{t}_j \leq t_j + 2\alpha$.
6. The subspace error $\text{SE}_t := \|(\mathbf{I} - \hat{\mathbf{P}}_t \hat{\mathbf{P}}_t') \mathbf{P}_t\|_2$ satisfies:

$$\text{SE}_t \leq \begin{cases} 1 & t \in [t_j, \hat{t}_j] \\ 10^{-2}\sqrt{\zeta} + 0.4^{k-1} & t \in [\hat{t}_j + (k-1)\alpha, \hat{t}_j + k\alpha - 1], \quad k = 1, 2, \dots, K \\ 10^{-2}\sqrt{\zeta} & t \in [\hat{t}_j + K\alpha, t_{j+1} - 1]. \end{cases}$$

3.8.3 Discussion

First we point out a significant limitation of the above result. Theorem 3.8.4 and Algorithm 4 assume that the clusters $\mathcal{G}_{j,k}$ are known a priori. This is not a very realistic assumption, as determining the clusters would require knowledge of all the eigenvalues of the true covariance matrix $\mathbf{\Lambda}_{a,t}$.

Let us compare the result for ReProCS (Algorithm 3) with that for ReProCS-cPCA (Algorithm 4) for the more general subspace change model. The ReProCS result needs $\kappa_{2s}([\mathbf{P}_0, \mathbf{P}_{1,\text{new}}, \dots, \mathbf{P}_{J,\text{new}}]) \leq 0.3$ while ReProCS-cPCA only needs $\max_j \kappa_{2s}(\mathbf{P}_{(j)}) \leq 0.3$

and $\max_j \kappa_{2s}(\mathbf{P}_{(j),\text{new}}) \leq 0.02$. Recall that $\mathbf{L} := [\ell_1, \ell_2, \dots, \ell_{t_{\max}}]$, $\mathbf{X} := [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{t_{\max}}]$, $r_{\text{mat}} := \text{rank}(\mathbf{L})$ and $s_{\text{mat}} = |\text{support}(\mathbf{X})|$. Also recall that in our model $s_{\text{mat}} \leq st_{\max}$.

Clearly $\text{rank}([\mathbf{P}_0, \mathbf{P}_{1,\text{new}}, \dots, \mathbf{P}_{J,\text{new}}]) = r_{\text{mat}}$, and so $r_{\text{mat}} \leq r_0 + Jc_{\text{new}}$. Thus if we assume $\kappa_1([\mathbf{P}_0, \mathbf{P}_{1,\text{new}}, \dots, \mathbf{P}_{J,\text{new}}])^2 \leq \mu r_{\text{mat}}/n$, and $\kappa_1(\mathbf{P}_{j,\text{new}}) \leq \tilde{\mu}c_{\text{new}}/n$ then ReProCS needs

$$2\mu s r_{\text{mat}} \leq 0.09n \text{ and } 2\tilde{\mu} s c_{\text{new}} \leq 0.0004n$$

where $r_{\text{mat}} \leq r_0 + Jc_{\text{new}}$. The support change model (Signal Model 3.2.11) requires $s \in \mathcal{O}(\frac{n}{\log n})$ and $J \in \mathcal{O}(n)$. Thus, if s grows as $n/\log n$, then r_{mat} can only grow as $\log n$. Thus, ReProCS allows

$$s_{\text{mat}} \in \mathcal{O}\left(\frac{nt_{\max}}{\log n}\right) \text{ and } r_{\text{mat}} \in \mathcal{O}(\log n)$$

As explained earlier this is a stronger requirement than PCP which allows

$$s_{\text{mat}} \in \mathcal{O}(nt_{\max}) \text{ and } r_{\text{mat}} \in \mathcal{O}\left(\frac{n}{(\log n)^2}\right)$$

On the other hand, using an argument similar to the one above, ReProCS-cPCA only needs

$$2\mu s(r_0 + c_{\text{dif}}) \leq 0.09n \text{ and } 2\tilde{\mu} s c_{\text{new}} \leq 0.0004n.$$

The support change model (Signal Model 3.2.11) requires $s \in \mathcal{O}(\frac{n}{\log n})$ and $J \in \mathcal{O}(n)$. Thus if s grows as $n/\log n$, r_0 can grow at most as $\log n$; however, c_{new} needs to be constant because of the lower bound on x_{\min} . Since $r_{\text{mat}} \leq r_0 + Jc_{\text{new}}$, this means r_{mat} can grow linearly with n . Thus, ReProCS-cPCA allows

$$s_{\text{mat}} \in \mathcal{O}\left(\frac{nt_{\max}}{\log n}\right) \text{ and } r_{\text{mat}} \in \mathcal{O}(n)$$

This requirement is comparable to what PCP needs: the requirement on s_{mat} is slightly stronger than PCP while that on r_{mat} is slightly weaker.

3.8.4 Cluster-PCA Algorithm (from [12])

In this section we present the ReProCS algorithm with cluster-PCA as Algorithm 4. It performs the same subspace change detection and estimation of new directions as ReProCS, but it differs in that at each subspace change time it performs a re-estimation of the previous subspace that effectively removes old directions from the estimate of \mathbf{P}_t .

One way to re-estimate the current subspace would be by standard PCA: at $t = \hat{t}_{j+1}$, compute

$$\hat{\mathbf{P}}_t \leftarrow \text{eigenvectors} \left(\sum_{\tau=\hat{t}_{j+1}-2\alpha-\tilde{\alpha}}^{\hat{t}_{j+1}-2\alpha} \hat{\ell}_\tau \hat{\ell}_\tau' \right)$$

Using the same proof method as used for Theorem 3.2.15, it can be shown that, as long as $f = \frac{\lambda^+}{\lambda^-}$ is small enough, we can show that this will give an accurate estimate of $\text{range}(\mathbf{P}_{(j)})$. However f cannot be small because our problem definition allows large noise, ℓ_t , but assumes slow subspace change. To resolve this problem, we recover the eigenvectors in clusters where the ratio of the largest to smallest eigenvalue in each cluster is small. We also have to assume that these clusters are sufficiently far apart from each other.

In a process that we call cluster-PCA, Algorithm 4 sequentially recovers eigenvectors corresponding to groups of eigenvalues that are close to each other and separated from the rest of the eigenvalues. The eigenvectors corresponding to the largest eigenvalues are recovered first. After this, the data is projected perpendicular to the estimated subspace to recover the eigenvectors corresponding to the next largest cluster. This process continues until all of $\text{span}(\mathbf{P}_{(j)})$ has been estimated. All of this takes place at $t = \hat{t}_{j+1}$ using the estimates $\hat{\ell}_t$ from the intervals $[(\hat{t}_{j+1} - \vartheta\tilde{\alpha} - 2\alpha) + (k-1)\tilde{\alpha} + 1, (\hat{t}_{j+1} - \vartheta\tilde{\alpha} - 2\alpha) + k\tilde{\alpha}]$, $k = 1, \dots, \vartheta$.

The motivation for the cluster-PCA step is Signal Model 3.2.8 where at the times when new directions are added to the subspace, some may also be removed. This is a more appropriate model when the sparse signal \mathbf{x}_t is the signal of interest. Under this model we are able to relax the requirements needed to prove a similar performance guarantee to Theorem 3.2.15. The difference is that instead of needing $\text{span}([\ell_1, \dots, \ell_{t_{\max}}])$ to be dense, we only require $\text{span}(\mathbf{P}_t)$ to be dense for all t . Of course when no directions are removed from the subspace, these are equivalent, because $\text{span}(\mathbf{P}_t) \subseteq \text{span}(\mathbf{P}_{t+1})$ for all t .

3.8.5 Proof of Theorem 3.8.4

Definition 3.8.5. *Redefine $\hat{\mathbf{P}}_{(j),*}$ to be the estimate of $\mathbf{P}_{(j),*} = \mathbf{P}_{(j-1)}$ after the cluster PCA step. That is $\hat{\mathbf{P}}_{(j),*} := \hat{\mathbf{P}}_{\hat{t}_j}$. Previously we had defined $\hat{\mathbf{P}}_{(j),*}$ to be the estimate of $\mathbf{P}_{(j),*}$ after the final projection PCA (addition) step, so this new definition is a natural extension. In both*

Algorithm 4 Recursive Projected CS with cluster-PCA (ReProCS-cPCA)

Parameters: algorithm parameters: $\xi, \omega, \alpha, \tilde{\alpha}, K$, model parameters: λ^-, ϑ , and $\tilde{c}_{j,k}$

Input: $n \times 1$ vector, \mathbf{m}_t , and $n \times r_0$ basis matrix $\hat{\mathbf{P}}_{(0)}$.

Output: $n \times 1$ vectors $\hat{\mathbf{x}}_t$ and $\hat{\boldsymbol{\ell}}_t$, a basis matrix $\hat{\mathbf{P}}_t, \hat{\mathbf{t}}_j, \hat{\mathbf{c}}_{j,\text{new},k}$.

In step 4(c)ii of Algorithm 3, include the following cluster-PCA (cPCA) step.

Re-estimate $\text{range}(\mathbf{P}_{(j-1)})$ by cluster-PCA (in order to remove the deleted directions)

1. set $\hat{\mathbf{G}}_{j-1,0} \leftarrow [\cdot]$
2. For $i = 1, 2, \dots, \vartheta$,
 - let $\hat{\mathbf{G}}_{j-1,\text{det},i} := [\hat{\mathbf{G}}_{j-1,1}, \hat{\mathbf{G}}_{j-1,2}, \dots, \hat{\mathbf{G}}_{j-1,i-1}]$ and compute

$$\mathcal{M}_{cpca} = (\mathbf{I} - \hat{\mathbf{G}}_{j-1,\text{det},i} \hat{\mathbf{G}}_{j-1,\text{det},i}') \left(\frac{1}{\tilde{\alpha}} \sum_{t \in \mathcal{I}_{j,k}} \hat{\boldsymbol{\ell}}_t \hat{\boldsymbol{\ell}}_t' \right) (\mathbf{I} - \hat{\mathbf{G}}_{j-1,\text{det},i} \hat{\mathbf{G}}_{j-1,\text{det},i}')$$

- $\hat{\mathbf{G}}_{j-1,i} \leftarrow \text{eigenvectors}(\mathcal{M}_{cpca}, \tilde{c}_{j,i})$

End for

3. set $\hat{\mathbf{P}}_t \leftarrow [\hat{\mathbf{G}}_{j-1,1} \cdots \hat{\mathbf{G}}_{j-1,\vartheta}]$, $\hat{\mathbf{P}}_{*,t} \leftarrow \hat{\mathbf{P}}_t$, and $\hat{\mathbf{P}}_{t,\text{new}} \leftarrow [\cdot]$.

The function $\text{eigenvectors}(\mathcal{M}, c)$ returns the matrix containing the eigenvectors corresponding to the c largest eigenvalues.

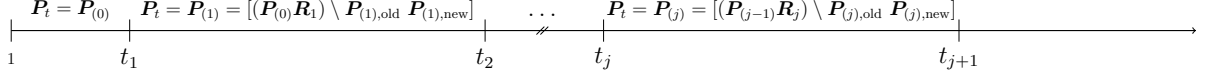


Figure 3.7 Signal Model 3.2.8

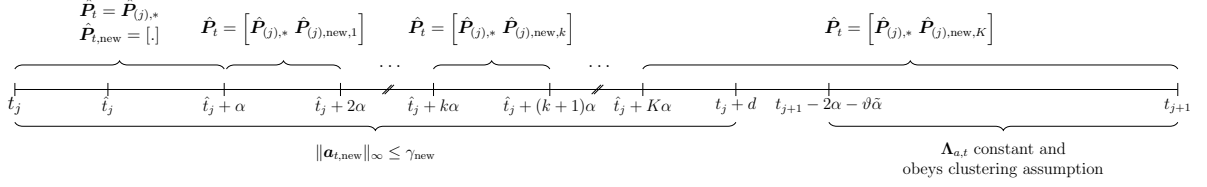


Figure 3.8 Algorithm 4

Figure 3.9: Diagrams for Signal Model 3.2.8 and Algorithm 4

cases, $\hat{P}_{(j),*}$ is the final estimate of $P_{(j-1)}$.

With the new definition of $\hat{P}_{(j),*}$, the statement of the definition of $\zeta_{j,*}$ remains the same:

$$\zeta_{j,*} := \text{dif}(\hat{P}_{(j),*}, P_{(j),*})$$

The tilde (\sim) will be used to indicate objects in the cluster PCA proof that have corresponding parts in the projection PCA proof.

Definition 3.8.6. *Redefine*

1. $\zeta_*^+ := r_{\max} \zeta$
2. $\tilde{\zeta}_k^+ := \frac{b_{\tilde{\mathcal{H}},k}}{b_{\tilde{A},k} - b_{\tilde{A},k,\perp} - b_{\tilde{\mathcal{H}},k}}$

Here $b_{\tilde{\mathbf{H}},k} := \tilde{b}_2 + \tilde{b}_{4,k} + \tilde{b}_{6,k}$. Where,

$$\begin{aligned}\tilde{b}_2 &:= \rho^2 h^+ (\phi^+)^2 (r_{\max} + c_{\text{new}})^2 \zeta^2 (b^2 r \gamma^2 + (1 - b^2) \lambda^+) + \epsilon \\ \tilde{b}_{4,k} &:= (r\zeta) \left(\frac{b^2}{\tilde{\alpha}(1 - b^2)} r \gamma^2 + (1 - b^2) \lambda^+ + \epsilon \right) (r_{\max} + c_{\text{new}}) \zeta \left(\sqrt{\rho^2 h^+} \phi^+ \right) + \\ &\quad \left(\frac{b^2}{\tilde{\alpha}(1 - b^2)} r \gamma^2 + (1 - b^2) \lambda_{j,k}^+ + \epsilon \right) (r_{\max} + c_{\text{new}}) \zeta \left(\sqrt{\rho^2 h^+} \phi^+ \right) \\ \tilde{b}_{6,k} &:= r\zeta \left(\frac{b^2}{\tilde{\alpha}(1 - b^2)} r \gamma^2 + 4\epsilon \right) + \frac{b^2}{\tilde{\alpha}(1 - b^2)} r \gamma^2 + 4\epsilon + (r\zeta)^2 \left(\frac{b^2}{\tilde{\alpha}(1 - b^2)} r \gamma^2 + (1 - b^2) \lambda_{j,k-1}^+ + \epsilon \right) + \\ &\quad r\zeta \left(\frac{b^2}{\tilde{\alpha}(1 - b^2)} r \gamma^2 + 4\epsilon \right) + \frac{(r\zeta)^3}{\sqrt{1 - r^2 \zeta^2}} \left(\frac{b^2}{\tilde{\alpha}(1 - b^2)} r \gamma^2 + 4\epsilon \right) + \\ &\quad \frac{(r\zeta)^2}{\sqrt{1 - r^2 \zeta^2}} \left(\frac{b^2}{\tilde{\alpha}(1 - b^2)} r \gamma^2 + (1 - b^2) \lambda_{j,k+1}^+ + \epsilon \right).\end{aligned}$$

Also, let

$$\begin{aligned}b_{\tilde{\mathbf{A}},k} &:= (1 - r^2 \zeta^2) [(1 - b^2) \lambda_k^- - \epsilon] - \sqrt{1 - r^2 \zeta^2} \left(\frac{b^2}{\tilde{\alpha}(1 - b^2)} r \gamma^2 + 4\epsilon \right) \left(r\zeta + \frac{r^2 \zeta^2}{\sqrt{1 - r^2 \zeta^2}} \right) \\ b_{\tilde{\mathbf{A}},k,\perp} &:= r^2 \zeta^2 \left(\frac{b^2}{\tilde{\alpha}(1 - b^2)} r \gamma^2 + (1 - b^2) \lambda_{j,k-1}^+ + \epsilon \right) - 2r\zeta \left(\frac{b^2}{\tilde{\alpha}(1 - b^2)} r \gamma^2 + (1 - b^2) \lambda_{j,k}^+ + \epsilon \right) - \\ &\quad \left(\frac{b^2}{\tilde{\alpha}(1 - b^2)} r \gamma^2 + (1 - b^2) \lambda_{j,k+1}^+ + \epsilon \right).\end{aligned}$$

To prove Theorem 3.8.4 we need to define a new event that says the $(j, k)^{\text{th}}$ cluster PCA step was successful.

Definition 3.8.7. Define

$$\text{CPCA}_{j,k}^a := \left\{ \left\| \left(\mathbf{I} - \sum_{i=1}^k \hat{\mathbf{G}}_{j,i} \hat{\mathbf{G}}_{j,i}' \right) \mathbf{G}_{j,k} \right\|_2 \leq \tilde{\zeta}_{j,k}^+ \right\}$$

and redefine

$$\Gamma_{j,0}^a := \Gamma_{j-1,\text{end}} \cap \text{DET}^a \cap \text{CPCA}_{j-1,1}^a \cap \dots \cap \text{CPCA}_{j-1,\vartheta}^a \quad \text{for } a = u_j \text{ or } a = u_j + 1$$

The definitions,

$$\Gamma_{j,k}^a := \Gamma_{j,k-1}^a \cap \text{PPCA}_{j,k}^a \quad \text{for } a = u_j \text{ or } a = u_j + 1$$

$$\Gamma_{j,\text{end}} := \left(\Gamma_{j,K}^{u_j} \cap \text{NODETS}_j^{u_j} \right) \cup \left(\Gamma_{j,K}^{u_j+1} \cap \text{NODETS}_j^{u_j+1} \right)$$

remain the same.

Lemma 3.8.8.

$$\mathbb{P}(\text{CPCA}_{j,k}^a \mid \Gamma_{j-1,\text{end}}, \text{DET}_j^a, \text{CPCA}_{j,1}^a, \dots, \text{CPCA}_{j,k-1}^a) \geq p_{\text{cpca}}$$

for $a = u_{j+1}$ or $a = u_{j+1} + 1$.

Fact 3.8.9.

1. If $\zeta_{j,*} \leq \zeta_{j,*}^+$ and $\zeta_{j,\text{new},k} \leq \zeta_{k,\text{new}}^+$, for $k = 1, \dots, K$ then $\text{dif}(\hat{\mathbf{P}}_{t_j+K\alpha}, \mathbf{P}_{(j)}) \leq \zeta_*^+ + c_{\text{new}}\zeta$
2. If $\tilde{\zeta}_{j,k} \leq \tilde{\zeta}_k^+$ for $k = 1, \dots, \vartheta$, then $\zeta_{j+1,*} := \text{dif}(\hat{\mathbf{P}}_{(j+1),*}, \mathbf{P}_{(j+1),*}) = \text{dif}(\hat{\mathbf{P}}_{t_{j+1}}, \mathbf{P}_{(j)}) \leq r_{\max}\zeta = \zeta_*^+$.
3. Thus the event $\Gamma_{j,\text{end}}$ implies $\text{dif}(\hat{\mathbf{P}}_{t_j+K\alpha}, \mathbf{P}_{(j)}) \leq \zeta_*^+ + c_{\text{new}}\zeta$. The event $\Gamma_{j,0}^a$ implies $\zeta_{j,*} \leq \zeta_*^+ = r_{\max}\zeta$ for $a = u_j$ or $a = u_{j+1}$. Thus the event $\Gamma_{j,k-1}^a$ also implies this.

Corollary 3.8.10. Combining Lemmas 3.5.11, 3.5.12, 3.5.14, and 3.8.8 gives

$$\begin{aligned} \mathbb{P}(\Gamma_{j,\text{end}} \mid \Gamma_{j-1,\text{end}}) &= \mathbb{P}\left(\left(\text{DET}^{u_j} \bigcap_{k=1}^{\vartheta} \text{CPCA}_{j-1,k}^{u_j} \bigcap_{k=1}^K \text{PPCA}_{j,k}^{u_j}\right) \cup \right. \\ &\quad \left. \left(\overline{\text{DET}^{u_j}} \cap \text{DET}^{u_j+1} \bigcap_{k=1}^{\vartheta} \text{CPCA}_{j-1,k}^{u_j+1} \bigcap_{k=1}^K \text{PPCA}_{j,k}^{u_j+1}\right) \mid \Gamma_{j-1,\text{end}}\right) \\ &= \mathbb{P}\left(\text{DET}^{u_j} \bigcap_{k=1}^{\vartheta} \text{CPCA}_{j-1,k}^{u_j} \bigcap_{k=1}^K \text{PPCA}_{j,k}^{u_j} \mid \Gamma_{j-1,\text{end}}\right) \\ &\quad + \mathbb{P}\left(\overline{\text{DET}^{u_j}} \cap \text{DET}^{u_j+1} \bigcap_{k=1}^{\vartheta} \text{CPCA}_{j-1,k}^{u_j+1} \bigcap_{k=1}^K \text{PPCA}_{j,k}^{u_j+1} \mid \Gamma_{j-1,\text{end}}\right) \\ &\geq p_{\text{det},0} \cdot (p_{\text{cpca}})^{\vartheta} \cdot (p_{\text{ppca}})^K + (1 - p_{\text{det},0}) \cdot p_{\text{det},1} \cdot (p_{\text{cpca}})^{\vartheta} \cdot (p_{\text{ppca}})^K \\ &\geq p_{\text{det},1} \cdot (p_{\text{cpca}})^{\vartheta} \cdot (p_{\text{ppca}})^K \end{aligned}$$

3.8.6 Proof of Lemma 3.8.8

Definition 3.8.11. Define

$$\tilde{\mathcal{I}}_{j,k} := \left[(\hat{t}_{j+1} - \vartheta\tilde{\alpha} - 2\alpha) + (k-1)\tilde{\alpha} + 1, (\hat{t}_{j+1} - \vartheta\tilde{\alpha} - 2\alpha) + k\tilde{\alpha} \right]$$

Definition 3.8.12. Define

$$\tilde{X}_{j,k} := [\boldsymbol{\nu}_1, \dots, \boldsymbol{\nu}_{(\hat{t}_{j+1} - \vartheta\tilde{\alpha} - 2\alpha) + k\tilde{\alpha}}]$$

and

$$\tilde{\Gamma}_{j,k}^{\hat{u}_j} := \Gamma_{j-1,\text{end}} \cap \text{DET}^{\hat{u}_j} \cap \text{CPACA}_{j-1,1}^{\hat{u}_j} \cap \cdots \cap \text{CPCA}_{j-1,k}^{\hat{u}_j}$$

for $\hat{u}_j = u_j$ or $\hat{u}_j = u_j + 1$.

Definition 3.8.13.

1. Let $\tilde{\mathbf{D}}_{j,k,\text{cur}} \stackrel{\text{QR}}{=} \tilde{\mathbf{E}}_{j,k,\text{cur}} \tilde{\mathbf{R}}_{j,k,\text{cur}}$ denote its reduced QR decomposition. So $\tilde{\mathbf{E}}_{j,k,\text{cur}}$ is a basis matrix, and $\tilde{\mathbf{R}}_{j,k,\text{cur}}$ is upper triangular. Let $\tilde{\mathbf{E}}_{j,k,\text{cur},\perp}$ be a basis matrix for the orthogonal complement of $\text{range}(\tilde{\mathbf{E}}_{j,k,\text{cur}})$.

2. Using $\tilde{\mathbf{E}}_{j,k,\text{cur}}$ and $\tilde{\mathbf{E}}_{j,k,\text{cur},\perp}$, define

$$\begin{aligned} \tilde{\mathbf{A}}_{j,k} &:= \frac{1}{\tilde{\alpha}} \sum_{t \in \tilde{\mathcal{I}}_{j,k}} \tilde{\mathbf{E}}_{j,k,\text{cur}}' \left(\mathbf{I} - \hat{\mathbf{G}}_{j,k,\text{det}} \hat{\mathbf{G}}_{j,k,\text{det}}' \right) \ell_t \ell_t' \left(\mathbf{I} - \hat{\mathbf{G}}_{j,k,\text{det}} \hat{\mathbf{G}}_{j,k,\text{det}}' \right) \tilde{\mathbf{E}}_{j,k,\text{cur}} \\ \tilde{\mathbf{A}}_{j,k,\perp} &:= \frac{1}{\tilde{\alpha}} \sum_{t \in \tilde{\mathcal{I}}_{j,k}} \tilde{\mathbf{E}}_{j,k,\text{cur},\perp}' \left(\mathbf{I} - \hat{\mathbf{G}}_{j,k,\text{det}} \hat{\mathbf{G}}_{j,k,\text{det}}' \right) \ell_t \ell_t' \left(\mathbf{I} - \hat{\mathbf{G}}_{j,k,\text{det}} \hat{\mathbf{G}}_{j,k,\text{det}}' \right) \tilde{\mathbf{E}}_{j,k,\text{cur},\perp} \end{aligned}$$

and let

$$\tilde{\mathbf{A}}_{j,k} := \begin{bmatrix} \tilde{\mathbf{E}}_{j,k,\text{cur}} & \tilde{\mathbf{E}}_{j,k,\text{cur},\perp} \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{A}}_{j,k} & \mathbf{0} \\ \mathbf{0} & \tilde{\mathbf{A}}_{j,k,\perp} \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{E}}_{j,k,\text{cur}}' \\ \tilde{\mathbf{E}}_{j,k,\text{cur},\perp}' \end{bmatrix}$$

3. Define

$$\tilde{\mathbf{H}}_{j,k} = \frac{1}{\tilde{\alpha}} \sum_{t \in \tilde{\mathcal{I}}_{j,k}} \left(\mathbf{I} - \hat{\mathbf{G}}_{j,k,\text{det}} \hat{\mathbf{G}}_{j,k,\text{det}}' \right) \hat{\ell}_t \hat{\ell}_t' \left(\mathbf{I} - \hat{\mathbf{G}}_{j,k,\text{det}} \hat{\mathbf{G}}_{j,k,\text{det}}' \right) - \tilde{\mathbf{A}}_{j,k}$$

From Algorithm 4,

$$\begin{aligned} &\frac{1}{\tilde{\alpha}} \sum_{t \in \tilde{\mathcal{I}}_{j,k}} \left(\mathbf{I} - \hat{\mathbf{G}}_{j,k,\text{det}} \hat{\mathbf{G}}_{j,k,\text{det}}' \right) \hat{\ell}_t \hat{\ell}_t' \left(\mathbf{I} - \hat{\mathbf{G}}_{j,k,\text{det}} \hat{\mathbf{G}}_{j,k,\text{det}}' \right) \\ &\stackrel{\text{EVD}}{=} \begin{bmatrix} \hat{\mathbf{G}}_{j,k} & \hat{\mathbf{G}}_{j,k,\perp} \end{bmatrix} \begin{bmatrix} \hat{\mathbf{\Lambda}}_{j,k} & \mathbf{0} \\ \mathbf{0} & \hat{\mathbf{\Lambda}}_{j,k,\perp} \end{bmatrix} \begin{bmatrix} \hat{\mathbf{G}}_{j,k}' \\ \hat{\mathbf{G}}_{j,k,\perp}' \end{bmatrix}. \end{aligned}$$

Define $\tilde{\zeta}_{j,k} := \text{dif}([\hat{\mathbf{G}}_{j,1}, \dots, \hat{\mathbf{G}}_{j,k}], \mathbf{G}_{j,k})$. Recall Lemma 3.5.17, which for the matrices above says that if $\lambda_{\min}(\tilde{\mathbf{A}}_{j,k}) - \|\tilde{\mathbf{A}}_{j,k,\perp}\|_2 - \|\tilde{\mathbf{H}}_{j,k}\|_2 > 0$, then

$$\tilde{\zeta}_{j,k} \leq \frac{\|\tilde{\mathbf{H}}_{j,k}\|_2}{\lambda_{\min}(\tilde{\mathbf{A}}_{j,k}) - \|\tilde{\mathbf{A}}_{j,k,\perp}\|_2 - \|\tilde{\mathbf{H}}_{j,k}\|_2} \quad (3.25)$$

Lemma 3.8.14.

$$\begin{aligned} \mathbb{P} \left(\lambda_{\min} \left(\tilde{\mathbf{A}}_{j,k} \right) \geq (1 - r^2 \zeta^2) \left[(1 - b^2) \lambda_k^- - \epsilon \right] \right. \\ \left. - \sqrt{1 - r^2 \zeta^2} \left(\frac{b^2}{\tilde{\alpha}(1 - b^2)} r \gamma^2 + 4\epsilon \right) \left(r \zeta + \frac{r^2 \zeta^2}{\sqrt{1 - r^2 \zeta^2}} \right) \mid \tilde{X}_{k-1} \right) \\ \geq 1 - rF(\tilde{\alpha}, \epsilon, c_{\text{new}} \gamma^2) - 2(3(2r)F(\tilde{\alpha}, \epsilon, 2r\gamma^2) - (2r)F(\tilde{\alpha}, \epsilon, 4r\gamma^2)). \end{aligned}$$

for all $\tilde{X}_{k-1} \in \tilde{\Gamma}_{k-1}$.

Lemma 3.8.15.

$$\begin{aligned} \mathbb{P} \left(\lambda_{\max} \left(\tilde{\mathbf{A}}_{j,k,\perp} \right) \leq r^2 \zeta^2 \left(\frac{b^2}{\tilde{\alpha}(1 - b^2)} r \gamma^2 + (1 - b^2) \lambda_{j,k-1}^+ + \epsilon \right) + \right. \\ \left. 2r\zeta \left(\frac{b^2}{\tilde{\alpha}(1 - b^2)} r \gamma^2 + (1 - b^2) \lambda_{j,k}^+ + \epsilon \right) + \right. \\ \left. \left(\frac{b^2}{\tilde{\alpha}(1 - b^2)} r \gamma^2 + (1 - b^2) \lambda_{j,k+1}^+ + \epsilon \right) \mid \tilde{X}_{k-1} \right) \geq 1 - 3rF(\tilde{\alpha}, \epsilon, r\gamma^2) \end{aligned}$$

for all $\tilde{X}_{k-1} \in \tilde{\Gamma}_{k-1}$.

Lemma 3.8.16.

$$\begin{aligned} \mathbb{P} \left(\|\tilde{\mathcal{H}}_{j,k}\|_2 \leq \tilde{b}_2 + 2\tilde{b}_{4,k} + 2\tilde{b}_{6,k} \mid \tilde{X}_{k-1} \right) \geq 1 - nF(\tilde{\alpha}, \epsilon, (\phi^+)^4 \zeta^2) \\ - rF(\tilde{\alpha}, \epsilon, r\gamma^2) \\ - 3 \left(3(2r)F(\tilde{\alpha}, \epsilon, 2r\gamma^2) - (2r)F(\tilde{\alpha}, \epsilon, 4r\gamma^2) \right) \\ - 2 \left(rF(\tilde{\alpha}, \epsilon, c_{\text{new}} \gamma^2) \right). \end{aligned}$$

for all $\tilde{X}_{k-1} \in \tilde{\Gamma}_{k-1}$.

Combining the above three lemmas gives Lemma 3.8.8.

3.9 Proofs of Lemmas 3.8.14, 3.8.15, and 3.8.16

3.9.1 Minor Lemmas for Proving the Main Lemmas

Definition 3.9.1. Define

1. $\mathbf{G}_{j,k,\text{det}} := [\mathbf{G}_{j,1}, \mathbf{G}_{j,2}, \dots, \mathbf{G}_{j,k-1}]$, $\hat{\mathbf{G}}_{j,k,\text{det}} := [\hat{\mathbf{G}}_{j,1}, \hat{\mathbf{G}}_{j,2}, \dots, \hat{\mathbf{G}}_{j,k-1}]$,
 $\mathbf{G}_{j,k,\text{undet}} := [\mathbf{G}_{j,k+1}, \mathbf{G}_{j,k+2}, \dots, \mathbf{G}_{j,\vartheta}]$, $\mathbf{G}_{j,k,\text{cur}} := \mathbf{G}_{j,k}$;
2. $\mathbf{a}_{t,\text{det}} := \mathbf{G}_{\text{det},k}' \ell_t$, $\mathbf{a}_{t,\text{cur}} := \mathbf{G}_{k,\text{cur}}' \ell_t$, and $\mathbf{a}_{t,\text{undet}} := \mathbf{G}_{\text{undet},k}' \ell_t$;
3. $\tilde{\mathbf{D}}_{j,k,\text{cur}} := (\mathbf{I} - \hat{\mathbf{G}}_{j,k,\text{det}} \hat{\mathbf{G}}_{j,k,\text{det}}') \mathbf{G}_{j,k}$, $\tilde{\mathbf{D}}_{j,k,\text{det}} := (\mathbf{I} - \hat{\mathbf{G}}_{j,k,\text{det}} \hat{\mathbf{G}}_{j,k,\text{det}}') \mathbf{G}_{j,k,\text{det}}$, $\tilde{\mathbf{D}}_{j,k,\text{undet}} := (\mathbf{I} - \hat{\mathbf{G}}_{j,k,\text{det}} \hat{\mathbf{G}}_{j,k,\text{det}}') \mathbf{G}_{j,k,\text{undet}}$;

Lemma 3.9.2 (Sparse Recovery for $\tilde{\mathcal{I}}_{j,k}$). Define $\Phi_{(j),K} := \mathbf{I} - \hat{\mathbf{P}}_{(j),*} \hat{\mathbf{P}}_{(j),*}' - \hat{\mathbf{P}}_{(j),\text{new},K} \hat{\mathbf{P}}_{(j),\text{new},K}'$. Assume that all conditions of Theorem 3.8.4 hold. If $\zeta_{j,*} \leq \zeta_{j,*}^+ := r_{\max} \zeta$ and $\zeta_{j,\text{new},K} \leq c\zeta$, then for all $t \in \tilde{\mathcal{I}}_{j,k}$,

1. $\|\Phi_{(j),K} \mathbf{P}_t\|_2 \leq (r_{\max} + c_{\text{new}}) \zeta$
2. $\|[(\Phi_{(j),K}) \tau_t' (\Phi_{(j),K}) \tau_t]^{-1}\|_2 \leq \phi^+ := 1.2$.
3. $\hat{\mathcal{T}}_t = \mathcal{T}_t$,
4. $\mathbf{e}_t = \mathbf{I}_{\mathcal{T}_t} [(\Phi_{(j),K}) \tau_t' (\Phi_{(j),K}) \tau_t]^{-1} \mathbf{I}_{\mathcal{T}_t}' \Phi_{(j),K} \mathbf{P}_t \mathbf{a}_t$ and $\|\mathbf{e}_t\|_2 \leq \phi^+ \sqrt{\zeta}$.

We will use the following fact throughout.

Fact 3.9.3. As before,

1. $\Gamma_{j-1,\text{end}}$ implies NODETS_{j-1}^a for $a = u_{j-1}$ or $a = u_{j-1} + 1$ which implies $\hat{t}_j \geq t_j$.
2. $\text{DET}_j^{u_j}$ implies that $\hat{t}_j \leq t_j + \alpha$.
3. $\text{DET}_j^{u_j+1}$ implies that $\hat{t}_j \leq t_j + 2\alpha$.
4. The above facts combined with the model assumption $d_2 \geq \vartheta \tilde{\alpha} + 2\alpha$ imply that conditioned on $\Gamma_{j-1,\text{end}} \cap \text{DET}_j^{u_j+1}$, $\tilde{\mathcal{I}}_{j,k} \subseteq [t_{j+1} - d_2, t_{j+1} - 1]$ for $j = 1, \dots, J$ and $k = 1, \dots, \vartheta$ and $\hat{u}_j = u_j$ or $\hat{u}_j = u_j + 1$. Thus, during these intervals, the covariance matrix $\mathbf{\Lambda}_{a,t}$ is assumed constant and obeys the clustering assumption.

Remark 3.9.4. As in the addition proof, we will remove the subscript j at various places.

The following two lemmas are proved the same way as their counterparts in Section 3.6 (Lemmas 3.6.2 and 3.6.3 respectively). As before, when they are applied, we will use $\epsilon = 0.01c_{\text{new}}\zeta\lambda^-$.

Lemma 3.9.5.

$$\begin{aligned}
1. \quad & \mathbb{P} \left(\lambda_{\max} \left(\frac{1}{\tilde{\alpha}} \sum_{t \in \tilde{\mathcal{I}}_{j,k}} \mathbf{a}_{t,\text{det}} \mathbf{a}_{t,\text{det}}' \right) \leq \frac{b^2}{\tilde{\alpha}(1-b^2)} r\gamma^2 + (1-b^2)\lambda_{j,k-1}^+ + \epsilon \mid \tilde{X}_{k-1} \right) \geq \\
& 1 - rF(\tilde{\alpha}, \epsilon, r\gamma^2) \\
2. \quad & \mathbb{P} \left(\lambda_{\max} \left(\frac{1}{\tilde{\alpha}} \sum_{t \in \tilde{\mathcal{I}}_{j,k}} \mathbf{a}_{t,\text{cur}} \mathbf{a}_{t,\text{cur}}' \right) \leq \frac{b^2}{\tilde{\alpha}(1-b^2)} r\gamma^2 + (1-b^2)\lambda_{j,k}^+ + \epsilon \mid \tilde{X}_{k-1} \right) \geq \\
& 1 - rF(\tilde{\alpha}, \epsilon, r\gamma^2) \\
3. \quad & \mathbb{P} \left(\lambda_{\max} \left(\frac{1}{\tilde{\alpha}} \sum_{t \in \tilde{\mathcal{I}}_{j,k}} \mathbf{a}_{t,\text{undet}} \mathbf{a}_{t,\text{undet}}' \right) \leq \frac{b^2}{\tilde{\alpha}(1-b^2)} r\gamma^2 + (1-b^2)\lambda_{j,k+1}^+ + \epsilon \mid \tilde{X}_{k-1} \right) \geq \\
& 1 - rF(\tilde{\alpha}, \epsilon, r\gamma^2) \\
4. \quad & \mathbb{P} \left(\lambda_{\min} \left(\frac{1}{\tilde{\alpha}} \sum_{t \in \tilde{\mathcal{I}}_{j,k}} \mathbf{a}_{t,\text{cur}} \mathbf{a}_{t,\text{cur}}' \right) \geq (1-b^2)\lambda_{j,k}^- - \epsilon \mid \tilde{X}_{k-1} \right) \geq \\
& 1 - rF(\tilde{\alpha}, \epsilon, c_{\text{new}}\gamma^2).
\end{aligned}$$

Results 1) and 2) also hold for $\mathbf{a}_{t,\text{det}}\mathbf{a}_t'$ and $\mathbf{a}_{t,\text{cur}}\mathbf{a}_t'$ respectively.

Lemma 3.9.6.

1.

$$\begin{aligned}
& \mathbb{P} \left(\left\| \frac{1}{\tilde{\alpha}} \sum_{t \in \tilde{\mathcal{I}}_{j,k}} \mathbf{a}_{t,\text{det}} \mathbf{a}_{t,\text{cur}}' \right\|_2 \leq \frac{b^2}{\tilde{\alpha}(1-b^2)} r\gamma^2 + 4\epsilon \mid \tilde{X}_{k-1} \right) \geq \\
& 1 - 3(2r)F(\tilde{\alpha}, \epsilon, 2r\gamma^2) - (2r)F(\tilde{\alpha}\epsilon, 4r\gamma^2).
\end{aligned}$$

2.

$$\begin{aligned}
& \mathbb{P} \left(\left\| \frac{1}{\tilde{\alpha}} \sum_{t \in \tilde{\mathcal{I}}_{j,k}} \mathbf{a}_{t,\text{det}} \mathbf{a}_{t,\text{undet}}' \right\|_2 \leq \frac{b^2}{\tilde{\alpha}(1-b^2)} r\gamma^2 + 4\epsilon \mid \tilde{X}_{k-1} \right) \geq \\
& 1 - 3(2r)F(\tilde{\alpha}, \epsilon, 2r\gamma^2) - (2r)F(\tilde{\alpha}\epsilon, 4r\gamma^2).
\end{aligned}$$

3.

$$\mathbb{P} \left(\left\| \frac{1}{\tilde{\alpha}} \sum_{t \in \tilde{\mathcal{I}}_{j,k}} \mathbf{a}_{t,\text{cur}} \mathbf{a}_{t,\text{undet}}' \right\|_2 \leq \frac{b^2}{\tilde{\alpha}(1-b^2)} r \gamma^2 + 4\epsilon \mid \tilde{X}_{k-1} \right) \geq 1 - 3(2r)F(\tilde{\alpha}, \epsilon, 2r\gamma^2) - (2r)F(\tilde{\alpha}\epsilon, 4r\gamma^2).$$

Lemma 3.9.7. [12] When $\tilde{X}_{j,k-1} \in \tilde{\Gamma}_{j,k-1}$,

1. $\|\tilde{\mathbf{D}}_{j,k,\text{det}}\|_2 \leq r_{\max}\zeta$
2. $\sqrt{1 - (r_{\max} + c_{\text{new}})^2\zeta^2} \leq \sigma_i(\tilde{\mathbf{R}}_{j,k}) \leq 1$
3. $\|\tilde{\mathbf{E}}_{j,k,\text{cur}}' \tilde{\mathbf{D}}_{j,k,\text{undet}}\|_2 \leq \frac{(r_{\max} + c_{\text{new}})^2\zeta^2}{\sqrt{1 - (r_{\max} + c_{\text{new}})^2\zeta^2}}$

We are now read to prove Lemmas 3.8.14, 3.8.15, and 3.8.16. Because the lemmas apply for all j , we will often omit the subscript j for convenience. Also, \sum_t will be used to mean $\sum_{t \in \tilde{\mathcal{I}}_{j,k}}$.

Proof of Lemma 3.8.14.

Recall $\tilde{\mathbf{A}}_k := \frac{1}{\tilde{\alpha}} \sum_{t \in \tilde{\mathcal{I}}_{j,k}} \tilde{\mathbf{E}}_k' \left(\mathbf{I} - \hat{\mathbf{G}}_{\text{det},k} \hat{\mathbf{G}}_{\text{det},k}' \right) \ell_t \ell_t' \left(\mathbf{I} - \hat{\mathbf{G}}_{\text{det},k} \hat{\mathbf{G}}_{\text{det},k}' \right) \tilde{\mathbf{E}}_k$. Observe that

$$\tilde{\mathbf{E}}_k' \left(\mathbf{I} - \sum_{i=1}^{k-1} \hat{\mathbf{G}}_{j,i} \hat{\mathbf{G}}_{j,i}' \right) \ell_t = \tilde{\mathbf{R}}_k \mathbf{a}_{t,\text{cur}} + \tilde{\mathbf{E}}_k' (\tilde{\mathbf{D}}_{\text{det},k} \mathbf{a}_{t,\text{det}} + \tilde{\mathbf{D}}_{\text{undet},k} \mathbf{a}_{t,\text{det}})$$

Let $\mathbf{Z}_t = \tilde{\mathbf{R}}_k \mathbf{a}_{t,\text{cur}}$ and $\mathbf{Y}_t = \tilde{\mathbf{E}}_k' (\tilde{\mathbf{D}}_{\text{det},k} \mathbf{a}_{t,\text{det}} + \tilde{\mathbf{D}}_{\text{undet},k} \mathbf{a}_{t,\text{undet}})$ Then

$$\tilde{\mathbf{A}}_k \succeq \frac{1}{\tilde{\alpha}} \sum_{t \in \tilde{\mathcal{I}}_{j,k}} \mathbf{Z}_t \mathbf{Z}_t' + \mathbf{Z}_t \mathbf{Y}_t' + \mathbf{Y}_t \mathbf{Z}_t' \quad (3.26)$$

Consider $\frac{1}{\tilde{\alpha}} \sum_t \mathbf{Z}_t \mathbf{Z}_t'$. Using a theorem of Ostrowski [27, Theorem 4.5.9], we have that

$$\lambda_{\min} \left(\frac{1}{\tilde{\alpha}} \sum_t \mathbf{Z}_t \mathbf{Z}_t' \right) \geq \lambda_{\min}(\tilde{\mathbf{R}}_k \tilde{\mathbf{R}}_k') \cdot \lambda_{\min} \left(\frac{1}{\tilde{\alpha}} \sum_t \mathbf{a}_{t,\text{cur}} \mathbf{a}_{t,\text{cur}}' \right)$$

By Lemmas 3.9.5 and 3.9.7, we get that

$$\mathbb{P} \left(\lambda_{\min} \left(\frac{1}{\tilde{\alpha}} \sum_t \mathbf{Z}_t \mathbf{Z}_t' \right) \geq (1 - r^2\zeta^2) [(1 - b^2)\lambda_k^- - \epsilon] \mid \tilde{X}_{k-1} \right) \geq 1 - rF(\tilde{\alpha}, \epsilon, c_{\text{new}}\gamma^2).$$

for all $\tilde{X}_{k-1} \in \tilde{\Gamma}_{k-1}$.

Next consider the term $\frac{1}{\tilde{\alpha}} \sum_t \mathbf{Z}_t \mathbf{Y}_t'$.

$$\begin{aligned} \frac{1}{\tilde{\alpha}} \sum_t \mathbf{Z}_t \mathbf{Y}_t' &= \frac{1}{\tilde{\alpha}} \sum_t \tilde{\mathbf{R}}_k \mathbf{a}_{t,\text{cur}} \left(\mathbf{a}_{t,\text{det}}' \tilde{\mathbf{D}}_{\text{det},k}' + \mathbf{a}_{t,\text{undet}}' \tilde{\mathbf{D}}_{\text{undet},k}' \right) \tilde{\mathbf{E}}_k \\ &= \frac{1}{\tilde{\alpha}} \sum_t \tilde{\mathbf{R}}_k \mathbf{a}_{t,\text{cur}} \mathbf{a}_{t,\text{det}}' \tilde{\mathbf{D}}_{\text{det},k}' \tilde{\mathbf{E}}_k + \tilde{\mathbf{R}}_k \mathbf{a}_{t,\text{cur}} \mathbf{a}_{t,\text{undet}}' \tilde{\mathbf{D}}_{\text{undet},k}' \tilde{\mathbf{E}}_k \end{aligned}$$

By Lemmas 3.9.6 and 3.9.7, we get that

$$\begin{aligned} \mathbb{P} \left(\lambda_{\min} \left(\frac{1}{\tilde{\alpha}} \sum_t \mathbf{Z}_t \mathbf{Y}_t' \right) \geq -\sqrt{1-r^2\zeta^2} \left(\frac{b^2}{\tilde{\alpha}(1-b^2)} r\gamma^2 + 4\epsilon \right) \left(r\zeta + \frac{r^2\zeta^2}{\sqrt{1-r^2\zeta^2}} \right) \mid \tilde{X}_{k-1} \right) \geq \\ 1 - 2(3(2r)F(\tilde{\alpha}, \epsilon, 2r\gamma^2) - (2r)F(\tilde{\alpha}\epsilon, 4r\gamma^2)). \end{aligned}$$

for all $\tilde{X}_{k-1} \in \tilde{\Gamma}_{k-1}$.

Therefore

$$\begin{aligned} \mathbb{P} \left(\lambda_{\min} \left(\tilde{\mathbf{A}}_{j,k} \right) \geq (1-r^2\zeta^2) [(1-b^2)\lambda_k^- - \epsilon] \right. \\ \left. - 2\sqrt{1-r^2\zeta^2} \left(\frac{b^2}{\tilde{\alpha}(1-b^2)} r\gamma^2 + 4\epsilon \right) \left(r\zeta + \frac{r^2\zeta^2}{\sqrt{1-r^2\zeta^2}} \right) \mid \tilde{X}_{k-1} \right) \\ \geq 1 - rF(\tilde{\alpha}, \epsilon, c_{\text{new}}\gamma^2) - 2(3(2r)F(\tilde{\alpha}, \epsilon, 2r\gamma^2) - (2r)F(\tilde{\alpha}\epsilon, 4r\gamma^2)). \end{aligned}$$

for all $\tilde{X}_{k-1} \in \tilde{\Gamma}_{k-1}$.

□

Proof of Lemma 3.8.15.

Recall $\tilde{\mathbf{A}}_{k,\perp} := \frac{1}{\tilde{\alpha}} \sum_{t \in \tilde{\mathcal{I}}_{j,k}} \tilde{\mathbf{E}}_{k,\perp}' \left(\mathbf{I} - \sum_{i=1}^{k-1} \hat{\mathbf{G}}_{j,i} \hat{\mathbf{G}}_{j,i}' \right) \ell_t \ell_t' \left(\mathbf{I} - \sum_{i=1}^{k-1} \hat{\mathbf{G}}_{j,i} \hat{\mathbf{G}}_{j,i}' \right) \tilde{\mathbf{E}}_{k,\perp}$. Notice that because $\tilde{\mathbf{E}}_{k,\perp}' \left(\mathbf{I} - \sum_{i=1}^{k-1} \hat{\mathbf{G}}_{j,i} \hat{\mathbf{G}}_{j,i}' \right) \mathbf{G}_k = \mathbf{0}$, $\tilde{\mathbf{E}}_{k,\perp}' \left(\mathbf{I} - \sum_{i=1}^{k-1} \hat{\mathbf{G}}_{j,i} \hat{\mathbf{G}}_{j,i}' \right) \ell_t = \tilde{\mathbf{E}}_{k,\perp}' (\tilde{\mathbf{D}}_{\text{det},k} \mathbf{a}_{t,\text{det}} + \tilde{\mathbf{D}}_{\text{undet},k} \mathbf{a}_{t,\text{undet}})$. So $\tilde{\mathbf{A}}_{k,\perp}$ can be written as

$$\tilde{\mathbf{A}}_{k,\perp} = \frac{1}{\tilde{\alpha}} \sum_t \tilde{\mathbf{E}}_{k,\perp}' (\tilde{\mathbf{D}}_{\text{det},k} \mathbf{a}_{t,\text{det}} + \tilde{\mathbf{D}}_{\text{undet},k} \mathbf{a}_{t,\text{undet}}) (\tilde{\mathbf{D}}_{\text{det},k} \mathbf{a}_{t,\text{det}} + \tilde{\mathbf{D}}_{\text{undet},k} \mathbf{a}_{t,\text{undet}})' \tilde{\mathbf{E}}_{k,\perp}.$$

Expanding the above gives

$$\begin{aligned} \tilde{\mathbf{A}}_{k,\perp} &= \frac{1}{\tilde{\alpha}} \sum_t \tilde{\mathbf{E}}_{k,\perp}' \left(\tilde{\mathbf{D}}_{\text{det},k} \mathbf{a}_{t,\text{det}} \mathbf{a}_{t,\text{det}}' \tilde{\mathbf{D}}_{\text{det},k}' + \tilde{\mathbf{D}}_{\text{det},k} \mathbf{a}_{t,\text{det}} \mathbf{a}_{t,\text{undet}}' \tilde{\mathbf{D}}_{\text{undet},k}' + \right. \\ &\quad \left. \tilde{\mathbf{D}}_{\text{undet},k} \mathbf{a}_{t,\text{undet}} \mathbf{a}_{t,\text{det}}' \tilde{\mathbf{D}}_{\text{det},k}' + \tilde{\mathbf{D}}_{\text{undet},k} \mathbf{a}_{t,\text{undet}} \mathbf{a}_{t,\text{undet}}' \tilde{\mathbf{D}}_{\text{undet},k}' \right) \tilde{\mathbf{E}}_{k,\perp}. \end{aligned}$$

By Lemmas 3.9.5, 3.9.6 and 3.9.7,

$$\begin{aligned} \mathbb{P} \left(\lambda_{\max} \left(\tilde{\mathbf{A}}_{j,k,\perp} \right) \leq r^2 \zeta^2 \left(\frac{b^2}{\tilde{\alpha}(1-b^2)} r \gamma^2 + (1-b^2) \lambda_{j,k-1}^+ + \epsilon \right) + \right. \\ \left. 2r\zeta \left(\frac{b^2}{\tilde{\alpha}(1-b^2)} r \gamma^2 + 4\epsilon \right) + \right. \\ \left. \left(\frac{b^2}{\tilde{\alpha}(1-b^2)} r \gamma^2 + (1-b^2) \lambda_{j,k+1}^+ + \epsilon \right) \mid \tilde{X}_{k-1} \right) \geq 1 - 3rF(\tilde{\alpha}, \epsilon, r\gamma^2) \end{aligned}$$

for all $\tilde{X}_{k-1} \in \tilde{\Gamma}_{k-1}$.

□

Proof of Lemma 3.8.16. For Ease of Notation, define $\Psi_{j,k} := (\mathbf{I} - \hat{\mathbf{G}}_{j,k,\det} \hat{\mathbf{G}}_{j,k,\det}')^{\dagger}$ and $\tilde{\mathbf{F}}_t = \tilde{\mathbf{E}}_k \tilde{\mathbf{E}}_k' \Psi_{k-1} \ell_t \ell_t' \Psi_{k-1} \tilde{\mathbf{E}}_{k,\perp} \tilde{\mathbf{E}}_{k,\perp}'$. Using the expression for $\tilde{\mathcal{H}}_{j,k}$ given in Definition 3.8.13, adding and subtracting $\tilde{\mathbf{F}}_t + \tilde{\mathbf{F}}_t'$, and noting that $\tilde{\mathbf{E}}_k \tilde{\mathbf{E}}_k' + \tilde{\mathbf{E}}_{k,\perp} \tilde{\mathbf{E}}_{k,\perp}' = \mathbf{I}$ we get that

$$\tilde{\mathcal{H}}_k = \frac{1}{\tilde{\alpha}} \sum_t \left(\Psi_{k-1} \mathbf{e}_t \mathbf{e}_t' \Psi_{k-1} - (\Psi_{k-1} \ell_t \mathbf{e}_t' \Psi_{k-1} + \Psi_{k-1} \mathbf{e}_t \ell_t' \Psi_{k-1}) + (\tilde{\mathbf{F}}_t + \tilde{\mathbf{F}}_t') \right)$$

Thus

$$\|\tilde{\mathcal{H}}_k\|_2 \leq \left\| \frac{1}{\tilde{\alpha}} \sum_t \mathbf{e}_t \mathbf{e}_t' \right\|_2 + 2 \left\| \frac{1}{\tilde{\alpha}} \sum_t \Psi_{k-1} \ell_t \mathbf{e}_t' \right\|_2 + 2 \left\| \frac{1}{\tilde{\alpha}} \sum_t \tilde{\mathbf{F}}_t \right\|_2 \quad (3.27)$$

1. First consider $\|\frac{1}{\tilde{\alpha}} \sum_t \mathbf{e}_t \mathbf{e}_t'\|_2$. By Lemma 3.9.2, when $\tilde{X}_{k-1} \in \tilde{\Gamma}_{k-1}$,

$$\mathbf{e}_t = \mathbf{I}_{\mathcal{T}_t} [(\Phi_{(j),K}) \mathcal{T}_t' (\Phi_{(j),K}) \mathcal{T}_t]^{-1} \mathbf{I}_{\mathcal{T}_t}' \Phi_{(j),K} \mathbf{P}_t \mathbf{a}_t$$

and $\|\mathbf{e}_t\|_2 \leq (\phi^+)^2 \zeta$. By the same argument as Lemma 3.6.1,

$$\begin{aligned} \mathbb{E}_{t-1} \left[\mathbf{e}_t \mathbf{e}_t' \mid \tilde{X}_{k-1} \right] &= \mathbf{I}_{\mathcal{T}_t} [(\Phi_{(j),K}) \mathcal{T}_t' (\Phi_{(j),K}) \mathcal{T}_t]^{-1} \mathbf{I}_{\mathcal{T}_t}' \Phi_{(j),K} \mathbf{P}_t (b^2 \mathbf{a}_{t-1} \mathbf{a}_{t-1}' + \Lambda_{\nu,t}) \\ &\quad \mathbf{P}_t' \Phi_{(j),K} \mathbf{I}_{\mathcal{T}_t} [(\Phi_{(j),K}) \mathcal{T}_t' (\Phi_{(j),K}) \mathcal{T}_t]^{-1} \mathbf{I}_{\mathcal{T}_t}' \\ &\preceq (\rho^2 h^+ (\phi^+)^2 (r_{\max} + c_{\text{new}})^2 \zeta^2 (b^2 r \gamma^2 + (1-b^2) \lambda^+)) \mathbf{I} \end{aligned}$$

Applying the Azuma corollary (Corollary 3.A.9) gives,

$$\begin{aligned} \mathbb{P} \left(\left\| \frac{1}{\tilde{\alpha}} \sum_t \mathbf{e}_t \mathbf{e}_t' \right\|_2 \leq \rho^2 h^+ (\phi^+)^2 (r_{\max} + c_{\text{new}})^2 \zeta^2 (b^2 r \gamma^2 + (1-b^2) \lambda^+) + \epsilon \mid \tilde{X}_{k-1} \right) \geq \\ 1 - nF(\tilde{\alpha}, \epsilon, (\phi^+)^4 \zeta^2) \end{aligned} \quad (3.28)$$

for all $\tilde{X}_{k-1} \in \tilde{\Gamma}_{k-1}$.

2. Next consider $\|\frac{1}{\tilde{\alpha}} \sum_t \Psi_{k-1} \ell_t \mathbf{e}_t'\|_2$.

$$\begin{aligned}
\Psi_{k-1} \ell_t \mathbf{e}_t' &= \left(\tilde{D}_k \mathbf{a}_{t,\text{cur}} + \tilde{D}_{\text{det},k} \mathbf{a}_{t,\text{det}} + \tilde{D}_{\text{undet},k} \mathbf{a}_{t,\text{undet}} \right) \mathbf{e}_t' \\
&= \left(\tilde{D}_{\text{det},k} \mathbf{a}_{t,\text{det}} + \begin{bmatrix} \tilde{D}_k & \tilde{D}_{\text{undet},k} \end{bmatrix} \begin{bmatrix} \mathbf{a}_{t,\text{cur}} \\ \mathbf{a}_{t,\text{undet}} \end{bmatrix} \right) \mathbf{e}_t' \\
&= \left(\tilde{D}_{\text{det},k} \mathbf{a}_{t,\text{det}} + \begin{bmatrix} \tilde{D}_k & \tilde{D}_{\text{undet},k} \end{bmatrix} \begin{bmatrix} \mathbf{a}_{t,\text{cur}} \\ \mathbf{a}_{t,\text{undet}} \end{bmatrix} \right) \\
&\quad (\mathbf{I}_{\mathcal{T}_t} [(\Phi_{(j),K})_{\mathcal{T}_t}' (\Phi_{(j),K})_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \Phi_{(j),K} \mathbf{P}_{(j)} \mathbf{a}_t)' \\
&= \left(\tilde{D}_{\text{det},k} \mathbf{a}_{t,\text{det}} \mathbf{a}_t' + \begin{bmatrix} \tilde{D}_k & \tilde{D}_{\text{undet},k} \end{bmatrix} \begin{bmatrix} \mathbf{a}_{t,\text{cur}} \\ \mathbf{a}_{t,\text{undet}} \end{bmatrix} \mathbf{a}_t' \right) \\
&\quad (\mathbf{I}_{\mathcal{T}_t} [(\Phi_{(j),K})_{\mathcal{T}_t}' (\Phi_{(j),K})_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}' \Phi_{(j),K} \mathbf{P}_{(j)})'
\end{aligned}$$

Let,

$$\mathbf{X}_t = \left(\tilde{D}_{\text{det},k} \mathbf{a}_{t,\text{det}} \mathbf{a}_t' + \begin{bmatrix} \tilde{D}_k & \tilde{D}_{\text{undet},k} \end{bmatrix} \begin{bmatrix} \mathbf{a}_{t,\text{cur}} \\ \mathbf{a}_{t,\text{undet}} \end{bmatrix} \mathbf{a}_t' \right) \mathbf{P}_{(j)}' \Phi_{(j),K}$$

and

$$\mathbf{Y}_t = \mathbf{I}_{\mathcal{T}_t} [(\Phi_{(j),K})_{\mathcal{T}_t}' (\Phi_{(j),K})_{\mathcal{T}_t}]^{-1} \mathbf{I}_{\mathcal{T}_t}'$$

Then by Lemma 3.A.3, which is stated and proved in the appendix,

$$\left\| \frac{1}{\tilde{\alpha}} \sum_t \Psi_{k-1} \ell_t \mathbf{e}_t' \right\|_2^2 \leq \lambda_{\max} \left(\frac{1}{\tilde{\alpha}} \sum_{t \in \tilde{\mathcal{I}}_{j,k}} \mathbf{X}_t \mathbf{X}_t' \right) \lambda_{\max} \left(\frac{1}{\tilde{\alpha}} \sum_{t \in \tilde{\mathcal{I}}_{j,k}} \mathbf{Y}_t \mathbf{Y}_t' \right)$$

The reason for using Lemma 3.A.3 is so that Lemma 3.3.3 can be applied to \mathbf{Y}_t . For \mathbf{X}_t ,

we use Lemmas 3.9.2, 3.9.5, and 3.9.7. This gives,

$$\begin{aligned}
\mathbb{P} \left(\left\| \frac{1}{\tilde{\alpha}} \sum_t \Psi_{k-1} \ell_t \mathbf{e}_t' \right\|_2 \leq (r\zeta) \left(\frac{b^2}{\tilde{\alpha}(1-b^2)} r\gamma^2 + (1-b^2)\lambda^+ + \epsilon \right) \right. \\
\left. (r_{\max} + c_{\text{new}}) \zeta \left(\sqrt{\rho^2 h^+} \phi^+ \right) + \right. \\
\left. \left(\frac{b^2}{\tilde{\alpha}(1-b^2)} r\gamma^2 + (1-b^2)\lambda_{j,k}^+ + \epsilon \right) (r_{\max} + c_{\text{new}}) \zeta \left(\sqrt{\rho^2 h^+} \phi^+ \right) \mid \tilde{X}_{k-1} \right) \\
\geq 1 - rF(\tilde{\alpha}, \epsilon, r\gamma^2)
\end{aligned}$$

for all $\tilde{X}_{k-1} \in \tilde{\Gamma}_{k-1}$.

3. Finally consider $\|\frac{1}{\tilde{\alpha}} \sum_t \tilde{\mathbf{F}}_t\|_2$. Recall that $\tilde{\mathbf{F}}_t = \tilde{\mathbf{E}}_k \tilde{\mathbf{E}}_k' \Psi_{k-1} \ell_t \ell_t' \Psi_{k-1} \tilde{\mathbf{E}}_{k,\perp} \tilde{\mathbf{E}}_{k,\perp}'$. We will use the fact that $\tilde{\mathbf{E}}_{k,\perp}' \tilde{\mathbf{D}}_k = \mathbf{0}$ to simplify this expression.

$$\begin{aligned} \tilde{\mathbf{F}}_t &= \tilde{\mathbf{E}}_k \tilde{\mathbf{E}}_k' \Psi_{k-1} \ell_t \ell_t' \Psi_{k-1} \tilde{\mathbf{E}}_{k,\perp} \tilde{\mathbf{E}}_{k,\perp}' \\ &= \tilde{\mathbf{E}}_k \tilde{\mathbf{E}}_k' (\tilde{\mathbf{D}}_k \mathbf{a}_{t,\text{cur}} + \tilde{\mathbf{D}}_{\text{det},k} \mathbf{a}_{t,\text{det}} + \tilde{\mathbf{D}}_{\text{undet},k} \mathbf{a}_{t,\text{undet}}) \\ &\quad (\tilde{\mathbf{D}}_{\text{det},k} \mathbf{a}_{t,\text{det}} + \tilde{\mathbf{D}}_{\text{undet},k} \mathbf{a}_{t,\text{undet}})' \tilde{\mathbf{E}}_{k,\perp} \tilde{\mathbf{E}}_{k,\perp}' \end{aligned}$$

Therefore,

$$\begin{aligned} \mathbb{P} \left(\left\| \frac{1}{\tilde{\alpha}} \sum_t \tilde{\mathbf{F}}_t \right\|_2 \leq r\zeta \left(\frac{b^2}{\tilde{\alpha}(1-b^2)} r\gamma^2 + 4\epsilon \right) + \right. \\ \left. \frac{b^2}{\tilde{\alpha}(1-b^2)} r\gamma^2 + 4\epsilon + \right. \\ \left. (r\zeta)^2 \left(\frac{b^2}{\tilde{\alpha}(1-b^2)} r\gamma^2 + (1-b^2)\lambda_{j,k-1}^+ + \epsilon \right) + \right. \\ \left. r\zeta \left(\frac{b^2}{\tilde{\alpha}(1-b^2)} r\gamma^2 + 4\epsilon \right) + \right. \\ \left. \frac{(r\zeta)^3}{\sqrt{1-r^2\zeta^2}} \left(\frac{b^2}{\tilde{\alpha}(1-b^2)} r\gamma^2 + 4\epsilon \right) + \right. \\ \left. \frac{(r\zeta)^2}{\sqrt{1-r^2\zeta^2}} \left(\frac{b^2}{\tilde{\alpha}(1-b^2)} r\gamma^2 + (1-b^2)\lambda_{j,k+1}^+ + \epsilon \right) \mid \tilde{X}_{k-1} \right) \geq \\ 1 - 3 \left(3(2r)F(\tilde{\alpha}, \epsilon, 2r\gamma^2) - (2r)F(\tilde{\alpha}\epsilon, 4r\gamma^2) \right) - 2 \left(rF(\tilde{\alpha}, \epsilon, c_{\text{new}}\gamma^2) \right) \end{aligned}$$

for all $\tilde{X}_{k-1} \in \tilde{\Gamma}_{k-1}$.

Therefore,

$$\begin{aligned} \mathbb{P} \left(\|\tilde{\mathcal{H}}_{j,k}\|_2 \leq \tilde{b}_2 + 2\tilde{b}_{4,k} + 2\tilde{b}_{6,k} \mid \tilde{X}_{k-1} \right) &\geq 1 - nF(\tilde{\alpha}, \epsilon, (\phi^+)^4\zeta^2) \\ &\quad - rF(\tilde{\alpha}, \epsilon, r\gamma^2) \\ &\quad - 3 \left(3(2r)F(\tilde{\alpha}, \epsilon, 2r\gamma^2) - (2r)F(\tilde{\alpha}\epsilon, 4r\gamma^2) \right) - 2 \left(rF(\tilde{\alpha}, \epsilon, c_{\text{new}}\gamma^2) \right). \end{aligned}$$

for all $\tilde{X}_{k-1} \in \tilde{\Gamma}_{k-1}$. The quantities \tilde{b}_2 , $\tilde{b}_{4,k}$, and $\tilde{b}_{6,k}$ were defined in Definition 3.8.6.

□

3.10 Simulations

In this section we show the results of simulation experiments that demonstrate the result we have proven. The results can be seen in Figures 3.10 and 3.11. In both cases, results are averaged over 50 simulations.

For both figures, we set $n = 256$ and $t_{\max} = 8200$. The low-dimensional vectors ℓ_t were generated the same way for both figures. There are 3 subspace changes, so $J = 3$. The subspace changed at times $t_1 = 701$, $t_2 = 3701$, and $t_3 = 6201$. The bounds on the entries of \mathbf{a}_t and $\mathbf{a}_{t,\text{new}}$ were $\gamma = 600$, and $\gamma_{\text{new}} = 5$ respectively. At each subspace change, 4 new directions were added, and 4 were removed. Therefore $c_{j,\text{new}} = c_{j,\text{old}} = 4$ for $j = 1, 2, 3$. The algorithm parameters were set as follows: $\alpha = 100$, $K = 12$, $\text{thresh} = \frac{1}{27}$, $\xi = \sqrt{2}$ and $\omega = 0.1$.

The difference in the two figures is how the supports of the sparse vectors \mathbf{x}_t are generated. In Figure 3.10, each index of \mathbf{x}_t was non-zero with probability 0.0586 (so the size of each support is around 15) independently of other indices and other times t . In Figure 3.11, all of the supports of \mathbf{x}_t were of size 10 and obeyed Signal Model 3.2.11 with $\varrho = 2$, and $\beta = 25$ (so $h^+ = \frac{25}{\alpha} = 0.25$).

For this simulated data, we compare ReProCS, ReProCS-cPCA, PCP, and mod-PCP. In order to compare with the batch methods PCP and mod-PCP, we ran those algorithms every 2α frames, using the observations \mathbf{m}_t from the previous 2α time instants. As one can see in Figure 3.10, the error made by ReProCS after a subspace change is very similar to that made by PCP. However, as ReProCS recovers the new directions, the error decays exponentially (notice that the y axis is logarithmic) as we have proven in Theorem 3.2.15.

In Figure 3.11, the supports of \mathbf{x}_t are highly correlated. This causes a problem for the batch methods and one can see that ReProCS has significantly better recovery compared to PCP and even mod-PCP when the supports of \mathbf{x}_t are correlated.

The last thing to notice is the difference between ReProCS and ReProCS-cPCA. Both of the figures demonstrate the results we have proven. That is, for both algorithms the error is initially large, but decays exponentially thereafter. Because ReProCS only adds new directions to its estimate of the subspace, the final error after a subspace change increases with j . The

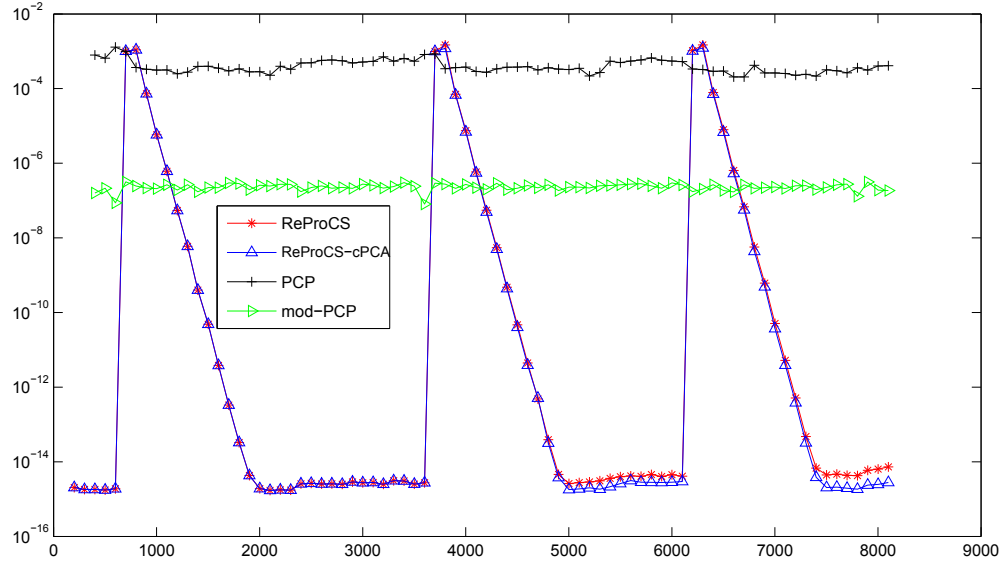


Figure 3.10: Support of \mathbf{X} determined by Bernoulli model. The y axis is $\frac{\|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2}{\|\mathbf{x}_t\|_2}$.

cluster PCA step in ReProCS-cPCA re-estimates the subspace, so that the error decays down to a value that does not increase with j .

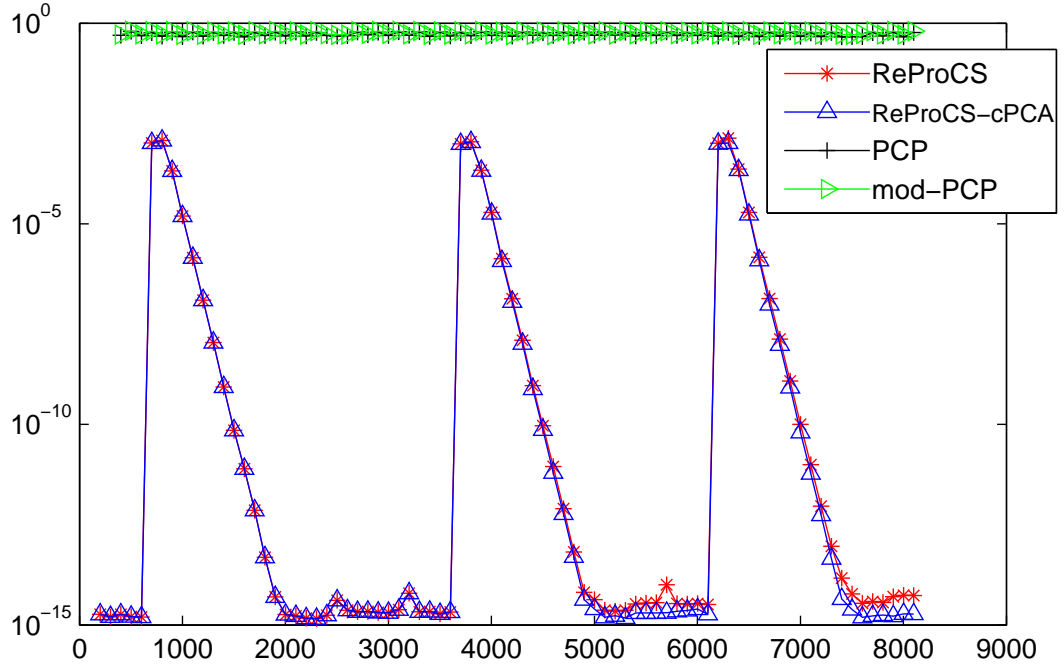


Figure 3.11: Support of \mathbf{X} obeys Signal Model 3.2.11. The y axis is $\frac{\|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2}{\|\mathbf{x}_t\|_2}$.

3.A Preliminaries

Lemma 3.A.1 (Exchanging the order of a double sum).

$$\sum_{t=0}^{\alpha-1} \sum_{i=0}^t x_t y_i = \sum_{i=0}^{\alpha-1} \sum_{t=i}^{\alpha-1} x_t y_i$$

Proof. Define [statement] to be the Boolean value of statement. Then,

$$\begin{aligned} \sum_{t=0}^{\alpha-1} \sum_{i=0}^t x_t y_i &= \sum_{t,i} [0 \leq i \leq t] [0 \leq t \leq \alpha-1] x_t y_i \\ &= \sum_{t,i} [0 \leq i \leq t \leq \alpha-1] x_t y_i \\ &= \sum_{t,i} [0 \leq i \leq \alpha-1] [i \leq t \leq \alpha-1] x_t y_i \\ &= \sum_{i=0}^{\alpha-1} \sum_{t=i}^{\alpha-1} x_t y_i \end{aligned}$$

□

Lemma 3.A.2 (Cauchy-Schwarz for a sum of vectors). *For vectors \mathbf{x}_t and \mathbf{y}_t ,*

$$\left(\sum_{t=1}^{\alpha} \mathbf{x}_t' \mathbf{y}_t \right)^2 \leq \left(\sum_t \|\mathbf{x}_t\|_2^2 \right) \left(\sum_t \|\mathbf{y}_t\|_2^2 \right)$$

Proof.

$$\begin{aligned} \left(\sum_{t=1}^{\alpha} \mathbf{x}_t' \mathbf{y}_t \right)^2 &= \left(\begin{bmatrix} \mathbf{x}_1' & \dots & \mathbf{x}_{\alpha}' \end{bmatrix} \begin{bmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_{\alpha} \end{bmatrix} \right)^2 \leq \left\| \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_{\alpha} \end{bmatrix} \right\|_2^2 \left\| \begin{bmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_{\alpha} \end{bmatrix} \right\|_2^2 = \\ &= \left(\sum_{t=1}^{\alpha} \|\mathbf{x}_t\|_2^2 \right) \left(\sum_{t=1}^{\alpha} \|\mathbf{y}_t\|_2^2 \right) \end{aligned}$$

The inequality is by Cauchy-Schwarz for a single vector. \square

Lemma 3.A.3 (Cauchy-Schwarz for a sum of matrices). *For matrices \mathbf{X}_t and \mathbf{Y}_t ,*

$$\left\| \frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbf{X}_t \mathbf{Y}_t' \right\|_2^2 \leq \lambda_{\max} \left(\frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbf{X}_t \mathbf{X}_t' \right) \lambda_{\max} \left(\frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbf{Y}_t \mathbf{Y}_t' \right)$$

Proof.

$$\begin{aligned} \left\| \sum_{t=1}^{\alpha} \mathbf{X}_t \mathbf{Y}_t' \right\|_2^2 &= \max_{\substack{\|\mathbf{x}\|=1 \\ \|\mathbf{y}\|=1}} \left| \mathbf{x}' \left(\sum_t \mathbf{X}_t \mathbf{Y}_t' \right) \mathbf{y} \right|^2 \\ &= \max_{\substack{\|\mathbf{x}\|=1 \\ \|\mathbf{y}\|=1}} \left| \sum_{t=1}^{\alpha} (\mathbf{X}_t' \mathbf{x})' (\mathbf{Y}_t' \mathbf{y}) \right|^2 \\ &\leq \max_{\substack{\|\mathbf{x}\|=1 \\ \|\mathbf{y}\|=1}} \left(\sum_{t=1}^{\alpha} \|\mathbf{X}_t' \mathbf{x}\|_2^2 \right) \left(\sum_{t=1}^{\alpha} \|\mathbf{Y}_t' \mathbf{y}\|_2^2 \right) \\ &= \max_{\|\mathbf{x}\|=1} \mathbf{x}' \sum_{t=1}^{\alpha} \mathbf{X}_t \mathbf{X}_t' \mathbf{x} \cdot \max_{\|\mathbf{y}\|=1} \mathbf{y}' \sum_{t=1}^{\alpha} \mathbf{Y}_t \mathbf{Y}_t' \mathbf{y} \\ &= \lambda_{\max} \left(\sum_{t=1}^{\alpha} \mathbf{X}_t \mathbf{X}_t' \right) \lambda_{\max} \left(\sum_{t=1}^{\alpha} \mathbf{Y}_t \mathbf{Y}_t' \right) \end{aligned}$$

The inequality is by Lemma 3.A.2. The penultimate line is because $\|\mathbf{v}\|_2^2 = \mathbf{v}' \mathbf{v}$ (used with $\mathbf{v} = \mathbf{X}_t' \mathbf{x}$). Multiplying both sides by $(\frac{1}{\alpha})^2$ gives the desired result. \square

Lemma 3.A.4. *Let X , Y , and Z be random variables. Assume that X is independent of $\{Y, Z\}$. Then*

$$\mathbb{E}[XY|Z] = \mathbb{E}[X]\mathbb{E}[Y|Z]$$

Proof. By the chain rule, $f_{X,Y|Z}(x,y|z) = f_{X|Y,Z}(x|y,z)f_{Y|Z}(y|z)$. Because X is independent of both Y and Z , $f_{X|Y,Z}(x|y,z) = f_X(x)$. \square

Remark 3.A.5. *Adopt the notation that if the lower limit of a sum exceeds the upper limit, then the sum is empty and therefore equal to zero.*

Lemma 3.A.6. *Let \mathbf{c}_τ be a sequence of vectors such that*

$$\mathbf{c}_\tau = b\mathbf{c}_{\tau-1} + \boldsymbol{\mu}_\tau$$

for a scalar b . Similarly, let $\tilde{\mathbf{c}}_\tau = b\tilde{\mathbf{c}}_{\tau-1} + \tilde{\boldsymbol{\mu}}_\tau$. Then

$$\sum_{\tau=0}^{\alpha-1} \mathbf{c}_\tau \tilde{\mathbf{c}}_\tau' = \sum_{i=0}^{\alpha-1} [\mathbf{Z}_{1,i} + \mathbf{Z}_{2,i} + \mathbf{Z}_{3,i} + \mathbf{Z}_{4,i}] + \mathbf{Z}_5$$

where

$$\begin{aligned} \mathbf{Z}_{1,i} &= \frac{(1 - b^{2(\alpha-i)})}{1 - b^2} \boldsymbol{\mu}_i \tilde{\boldsymbol{\mu}}_i', \\ \mathbf{Z}_{2,i} &= \sum_{i_2=0}^{i-1} \frac{(1 - b^{2(\alpha-i)})}{1 - b^2} b^{i-i_2} \boldsymbol{\mu}_i \tilde{\boldsymbol{\mu}}_{i_2}', \\ \mathbf{Z}_{3,i} &= \sum_{i_2=\alpha-i}^{\alpha-1} \frac{(1 - b^{2(\alpha-i_2)})}{1 - b^2} b^{i+i_2-\alpha+1} \boldsymbol{\mu}_{\alpha-i-1} \tilde{\boldsymbol{\mu}}_{i_2}', \\ \mathbf{Z}_{4,i} &= \frac{b^{i+1}(1 - b^{2(\alpha-i)})}{1 - b^2} (\boldsymbol{\mu}_i \tilde{\mathbf{c}}_{-1}' + \mathbf{c}_{-1} \tilde{\boldsymbol{\mu}}_i') \\ \mathbf{Z}_5 &= \frac{b^2(1 - b^{2\alpha})}{1 - b^2} \mathbf{c}_{-1} \tilde{\mathbf{c}}_{-1}' \end{aligned}$$

Proof. We start with some simple expansions.

$$\begin{aligned} \sum_{\tau=0}^{\alpha-1} \mathbf{c}_\tau \tilde{\mathbf{c}}_\tau' &= \sum_{\tau=0}^{\alpha-1} (b\mathbf{c}_{\tau-1} + \boldsymbol{\mu}_\tau)(b\tilde{\mathbf{c}}_{\tau-1} + \tilde{\boldsymbol{\mu}}_\tau)' \\ &= \sum_{\tau=0}^{\alpha-1} \left(b^{\tau+1} \mathbf{c}_{-1} + \sum_{i=0}^{\tau} b^{\tau-i} \boldsymbol{\mu}_i \right) \left(b^{\tau+1} \tilde{\mathbf{c}}_{-1} + \sum_{i=0}^{\tau} b^{\tau-i} \tilde{\boldsymbol{\mu}}_i \right)' \\ &= \sum_{\tau=0}^{\alpha-1} b^{2(\tau+1)} \mathbf{c}_{-1} \tilde{\mathbf{c}}_{-1}' + \sum_{\tau=0}^{\alpha-1} \sum_{i=0}^{\tau} \sum_{i_2=0}^{\tau} b^{2\tau-i-i_2} \boldsymbol{\mu}_i \tilde{\boldsymbol{\mu}}_{i_2}' + \\ &\quad \sum_{\tau=0}^{\alpha-1} \sum_{i=0}^{\tau} b^{2\tau+1-i} \mathbf{c}_{-1} \tilde{\boldsymbol{\mu}}_i' + b^{2\tau+1-i} \boldsymbol{\mu}_i \tilde{\mathbf{c}}_{-1}' \end{aligned} \tag{3.29}$$

By summing over the b 's, the first term is equal to \mathbf{Z}_5 . Applying Lemma 3.A.1 to the cross terms and summing over the b 's, we get

$$\begin{aligned} \sum_{\tau=0}^{\alpha-1} \sum_{i=0}^{\tau} b^{2\tau+1-i} \mathbf{c}_{-1} \tilde{\boldsymbol{\mu}}_i' + b^{2\tau+1-i} \boldsymbol{\mu}_i \tilde{\mathbf{c}}_{-1}' &= \sum_{i=0}^{\alpha-1} \sum_{\tau=i}^{\alpha-1} b^{2\tau+1-i} (\mathbf{c}_{-1} \tilde{\boldsymbol{\mu}}_i' + \boldsymbol{\mu}_i \tilde{\mathbf{c}}_{-1}') \\ &= \sum_{i=0}^{\alpha-1} \frac{b^{i+1}(1-b^{2(\alpha-i)})}{1-b^2} (\boldsymbol{\mu}_i \tilde{\mathbf{c}}_{-1}' + \mathbf{c}_{-1} \tilde{\boldsymbol{\mu}}_i') \\ &= \sum_{i=0}^{\alpha-1} \mathbf{Z}_{4,i} \end{aligned}$$

The remaining three terms, $\mathbf{Z}_{1,i}$, $\mathbf{Z}_{2,i}$, and $\mathbf{Z}_{3,i}$ will all come from the middle term in (3.29)

We apply Lemma 3.A.1 to term2, split the innermost sum into 3 parts, and use the fact that

$$\sum_{i=0}^{\alpha-1} x_i = \sum_{i=0}^{\alpha-1} x_{\alpha-1-i}$$

$$\begin{aligned} \sum_{\tau=0}^{\alpha-1} \sum_{i=0}^{\tau} \sum_{i_2=0}^{\tau} b^{2\tau-i-i_2} \boldsymbol{\mu}_i \tilde{\boldsymbol{\mu}}_{i_2}' &= \sum_{i=0}^{\alpha-1} \sum_{\tau=i}^{\alpha-1} \sum_{i_2=0}^{\tau} b^{2\tau-i-i_2} \boldsymbol{\mu}_i \tilde{\boldsymbol{\mu}}_{i_2}' \\ \text{(split inner sum)} &= \sum_{i=0}^{\alpha-1} \sum_{\tau=i}^{\alpha-1} \sum_{i_2=0}^{i-1} b^{2\tau-i-i_2} \boldsymbol{\mu}_i \tilde{\boldsymbol{\mu}}_{i_2}' + \sum_{i=0}^{\alpha-1} \sum_{\tau=i}^{\alpha-1} b^{2\tau-2i} \boldsymbol{\mu}_i \tilde{\boldsymbol{\mu}}_i' + \\ &\quad \sum_{i=0}^{\alpha-1} \sum_{\tau=i}^{\alpha-1} \sum_{i_2=i+1}^{\tau} b^{2\tau-i-i_2} \boldsymbol{\mu}_i \tilde{\boldsymbol{\mu}}_{i_2}' \\ \text{(switch indices)} &= \sum_{i=0}^{\alpha-1} \sum_{i_2=0}^{i-1} \sum_{\tau=i}^{\alpha-1} b^{2\tau-i-i_2} \boldsymbol{\mu}_i \tilde{\boldsymbol{\mu}}_{i_2}' + \sum_{i=0}^{\alpha-1} \sum_{\tau=i}^{\alpha-1} b^{2\tau-2i} \boldsymbol{\mu}_i \tilde{\boldsymbol{\mu}}_i' + \\ &\quad \sum_{i=0}^{\alpha-1} \sum_{i_2=i+1}^{\alpha-1} \sum_{\tau=i_2}^{\alpha-1} b^{2\tau-i-i_2} \boldsymbol{\mu}_i \tilde{\boldsymbol{\mu}}_{i_2}' \\ \text{(apply fact)} &= \sum_{i=0}^{\alpha-1} \sum_{i_2=0}^{i-1} \sum_{\tau=i}^{\alpha-1} b^{2\tau-i-i_2} \boldsymbol{\mu}_i \tilde{\boldsymbol{\mu}}_{i_2}' + \sum_{i=0}^{\alpha-1} \sum_{\tau=i}^{\alpha-1} b^{2\tau-2i} \boldsymbol{\mu}_i \tilde{\boldsymbol{\mu}}_i' \\ &\quad + \sum_{i=0}^{\alpha-1} \sum_{i_2=\alpha-i}^{\alpha-1} \sum_{\tau=i_2}^{\alpha-1} b^{2\tau-(\alpha-1-i)-i_2} \boldsymbol{\mu}_{\alpha-1-i} \tilde{\boldsymbol{\mu}}_{i_2}' \\ \text{(geometric series)} &= \sum_{i=0}^{\alpha-1} [\mathbf{Z}_{2,i} + \mathbf{Z}_{1,i} + \mathbf{Z}_{3,i}] \end{aligned}$$

□

Fact 3.A.7. For an event \mathcal{E} and random variable X , $\mathbb{P}(\mathcal{E}|X) \geq p$ for all $X \in \mathcal{C}$ implies that

$$\mathbb{P}(\mathcal{E}|X \in \mathcal{C}) \geq p.$$

Theorem 3.A.8 (Matrix Azuma). *[11, Theorem 7.1] Consider a finite adapted sequence \mathbf{Z}_t of $n \times n$ random Hermitian matrices, and a fixed sequence \mathbf{A}_t of Hermitian matrices that satisfy*

$$\mathbb{E}_{t-1}[\mathbf{Z}_t] = \mathbf{0} \quad \text{and} \quad \mathbf{Z}_t^2 \preceq \mathbf{A}_t^2 \quad \text{with probability 1.}$$

Define the variance parameter

$$\sigma^2 := \left\| \sum_t \mathbf{A}_t^2 \right\|_2.$$

Then, for all $\epsilon > 0$,

$$\mathbb{P} \left(\lambda_{\max} \left(\sum_t \mathbf{Z}_t \right) \geq \epsilon \right) \leq n \exp \left(\frac{-\epsilon^2}{8\sigma^2} \right)$$

The following corollary extends the above result to the case where $\mathbb{E}_{t-1}[\mathbf{Z}_t] \neq \mathbf{0}$ and also includes conditioning on another random variable.

Corollary 3.A.9 (Matrix Azuma conditioned on another random variable for a nonzero mean Hermitian matrix). *Consider an α -length sequence $\{\mathbf{Z}_t\}_{1 \leq t \leq \alpha}$ of random Hermitian matrices of size $n \times n$ given a random variable X . Assume that, for all $X \in \mathcal{C}$, (i) $\mathbb{P}(b_1 \mathbf{I} \preceq \mathbf{Z}_t \preceq b_2 \mathbf{I} | X) = 1$, for $1 \leq t \leq \alpha$ and (ii) $b_3 \mathbf{I} \preceq \frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbb{E}_{t-1}(\mathbf{Z}_t | X) \preceq b_4 \mathbf{I}$. Then for all $\epsilon > 0$,*

$$\begin{aligned} \mathbb{P} \left(\lambda_{\max} \left(\frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbf{Z}_t \right) \leq b_4 + \epsilon \middle| X \right) &\geq 1 - n \exp \left(\frac{-\alpha \epsilon^2}{8(b_2 - b_1)^2} \right) \\ \mathbb{P} \left(\lambda_{\min} \left(\frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbf{Z}_t \right) \geq b_3 - \epsilon \middle| X \right) &\geq 1 - n \exp \left(\frac{-\alpha \epsilon^2}{8(b_2 - b_1)^2} \right) \end{aligned}$$

Proof. 1. Let $\mathbf{Y}_t := \mathbf{Z}_t - \mathbb{E}_{t-1}(\mathbf{Z}_t | X)$. Clearly $\mathbb{E}_{t-1}(\mathbf{Y}_t | X) = \mathbf{0}$. Since for all $X \in \mathcal{C}$, $\mathbb{P}(b_1 \mathbf{I} \preceq \mathbf{Z}_t \preceq b_2 \mathbf{I} | X) = 1$ and since for an Hermitian matrix, $\lambda_{\max}(\cdot)$ is a convex function, and $\lambda_{\min}(\cdot)$ is a concave function, $b_1 \mathbf{I} \preceq \mathbb{E}_{t-1}(\mathbf{Z}_t | X) \preceq b_2 \mathbf{I}$ for all $X \in \mathcal{C}$. Therefore, $\mathbb{P}(\mathbf{Y}_t^2 \preceq (b_2 - b_1)^2 \mathbf{I} | X) = 1$ for all $X \in \mathcal{C}$. Thus, for Theorem 3.A.8, $\sigma^2 = \left\| \sum_{t=1}^{\alpha} (b_2 - b_1)^2 \mathbf{I} \right\|_2 = \alpha(b_2 - b_1)^2$. For any $X \in \mathcal{C}$, applying Theorem 3.A.8 for $\{\mathbf{Y}_t\}_{t=1, \dots, \alpha}$ conditioned on X , we get that, for any $\epsilon > 0$,

$$\mathbb{P} \left(\lambda_{\max} \left(\frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbf{Y}_t \right) \leq \epsilon \middle| X \right) > 1 - n \exp \left(\frac{-\alpha \epsilon^2}{8(b_2 - b_1)^2} \right) \quad \text{for all } X \in \mathcal{C}$$

By Weyl's theorem, $\lambda_{\max}(\frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbf{Y}_t) = \lambda_{\max}(\frac{1}{\alpha} \sum_{t=1}^{\alpha} (\mathbf{Z}_t - \mathbb{E}_{t-1}(\mathbf{Z}_t | X))) \geq \lambda_{\max}(\frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbf{Z}_t) + \lambda_{\min}(\frac{1}{\alpha} \sum_{t=1}^{\alpha} -\mathbb{E}_{t-1}(\mathbf{Z}_t | X))$.

Since $\lambda_{\min}(\frac{1}{\alpha} \sum_{t=1}^{\alpha} -\mathbb{E}_{t-1}(\mathbf{Z}_t|X)) = -\lambda_{\max}(\frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbb{E}_{t-1}(\mathbf{Z}_t|X)) \geq -b_4$,
 thus $\lambda_{\max}(\frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbf{Y}_t) \geq \lambda_{\max}(\frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbf{Z}_t) - b_4$. Therefore,

$$\mathbb{P}\left(\lambda_{\max}\left(\frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbf{Z}_t\right) \leq b_4 + \epsilon \middle| X\right) > 1 - n \exp\left(\frac{-\alpha\epsilon^2}{8(b_2 - b_1)^2}\right) \text{ for all } X \in \mathcal{C}$$

2. Now let $\mathbf{Y}_t = \mathbb{E}_{t-1}(\mathbf{Z}_t|X) - \mathbf{Z}_t$. As before, $\mathbb{E}_{t-1}(\mathbf{Y}_t|X) = 0$ and conditioned on any $X \in \mathcal{C}$, $\mathbf{P}(Y_t^2 \preceq (b_2 - b_1)^2 \mathbf{I} | X) = 1$. As before, applying Theorem 3.A.8, we get that for any $\epsilon > 0$,

$$\mathbb{P}\left(\lambda_{\max}\left(\frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbf{Y}_t\right) \leq \epsilon \middle| X\right) > 1 - n \exp\left(\frac{-\alpha\epsilon^2}{8(b_2 - b_1)^2}\right) \text{ for all } X \in \mathcal{C}$$

By Weyl's theorem, $\lambda_{\max}(\frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbf{Y}_t) = \lambda_{\max}(\frac{1}{\alpha} \sum_{t=1}^{\alpha} (\mathbb{E}_{t-1}(\mathbf{Z}_t|X) - \mathbf{Z}_t))$
 $\geq \lambda_{\min}(\frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbb{E}_{t-1}(\mathbf{Z}_t|X)) + \lambda_{\max}(\frac{1}{\alpha} \sum_{t=1}^{\alpha} -\mathbf{Z}_t) =$
 $\lambda_{\min}(\frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbb{E}_{t-1}(\mathbf{Z}_t|X)) - \lambda_{\min}(\frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbf{Z}_t) \geq b_3 - \lambda_{\min}(\frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbf{Z}_t)$ Therefore, for any $\epsilon > 0$,

$$\mathbb{P}\left(\lambda_{\min}\left(\frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbf{Z}_t\right) \geq b_3 - \epsilon \middle| X\right) \geq 1 - n \exp\left(\frac{-\alpha\epsilon^2}{8(b_2 - b_1)^2}\right) \text{ for all } X \in \mathcal{C}$$

□

We can further extend this to the case of a matrix which is not necessarily Hermitian.

Corollary 3.A.10 (Matrix Azuma conditioned on another random variable for an arbitrary matrix). *Consider an α -length adapted sequence $\{\mathbf{Z}_t\}$ of random matrices of size $n_1 \times n_2$ given a random variable X . Assume that, for all $X \in \mathcal{C}$, (i) $\mathbb{P}(\|\mathbf{Z}_t\|_2 \leq b_1 | X) = 1$ and (ii) $\|\frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbb{E}_{t-1}(\mathbf{Z}_t|X)\|_2 \leq b_2$. Then, for all $\epsilon > 0$,*

$$\mathbb{P}\left(\left\|\frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbf{Z}_t\right\|_2 \leq b_2 + \epsilon \middle| X\right) \geq 1 - (n_1 + n_2) \exp\left(\frac{-\alpha\epsilon^2}{8(2b_1)^2}\right)$$

Proof. Define the dilation of an $n_1 \times n_2$ matrix \mathbf{M} as $\text{dilation}(\mathbf{M}) := \begin{bmatrix} \mathbf{0} & \mathbf{M}' \\ \mathbf{M} & \mathbf{0} \end{bmatrix}$. Notice that this is an $(n_1 + n_2) \times (n_1 + n_2)$ Hermitian matrix [11]. As shown in [11, equation 2.12],

$$\lambda_{\max}(\text{dilation}(\mathbf{M})) = \|\text{dilation}(\mathbf{M})\|_2 = \|\mathbf{M}\|_2 \quad (3.30)$$

Thus, the corollary assumptions imply that $\mathbf{P}(\|\text{dilation}(\mathbf{Z}_t)\|_2 \leq b_1 | X) = 1$ for all $X \in \mathcal{C}$. By (3.30) and the definition of dilation,

$$\frac{1}{\alpha} \sum_t \mathbb{E}_{t-1}[\text{dilation}(\mathbf{Z}_t) | X] = \text{dilation} \left(\frac{1}{\alpha} \sum_t \mathbb{E}_{t-1}[\mathbf{Z}_t | X] \right) \preceq b_2 \mathbf{I}$$

Thus, applying Corollary 3.A.9 to the sequence $\{\text{dilation}(\mathbf{Z}_t)\}_{t=1, \dots, \alpha}$, we get that,

$$\mathbb{P} \left(\lambda_{\max} \left(\frac{1}{\alpha} \sum_{t=1}^{\alpha} \text{dilation}(\mathbf{Z}_t) \right) \leq b_2 + \epsilon \middle| X \right) \geq 1 - (n_1 + n_2) \exp \left(\frac{-\alpha \epsilon^2}{32 b_1^2} \right) \text{ for all } X \in \mathcal{C}$$

Using (3.30), $\lambda_{\max}(\frac{1}{\alpha} \sum_{t=1}^{\alpha} \text{dilation}(\mathbf{Z}_t)) = \lambda_{\max}(\text{dilation}(\frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbf{Z}_t)) = \|\frac{1}{\alpha} \sum_{t=1}^{\alpha} \mathbf{Z}_t\|_2$ gives the final result. \square

3.B Proofs of Support Change Lemmas

Proof of Lemma 3.2.18. We now provide the complete proof for 1) and 2) in the proof sketch.

For 1) we have by the bound in [24]

$$\begin{aligned}
& \mathbb{P}(\text{The object moves at least once every } \beta \text{ instants in the interval } \mathcal{J}_u) \\
&= \mathbb{P}(\text{The bit sequence } \theta_{(u-1)\alpha+1} \dots \theta_{u\alpha} \text{ does not contain a sequence of } \beta \text{ consecutive zeros}) \\
&\geq \left(1 - (1 - q)^\beta\right)^{\alpha - \beta + 1} \\
&\geq \left(1 - (1 - q)^\beta\right)^\alpha.
\end{aligned}$$

We need the object to move at least once every β time instants in *every* interval \mathcal{J}_u . We have

$$\begin{aligned}
\mathbb{P}(\text{The object moves at least once every } \beta \text{ instants in every } \mathcal{J}_u) &\geq \left(1 - (1 - q)^\beta\right)^{\lceil \frac{t_{\max}}{\alpha} \rceil \alpha} \\
&\geq \left(1 - (1 - q)^\beta\right)^{(\frac{t_{\max}}{\alpha} + 1)\alpha} \\
&\geq 1 - (t_{\max} + \alpha)(1 - q)^\beta.
\end{aligned}$$

This probability will be greater than $1 - \frac{n^{-10}}{2}$ if

$$q \geq 1 - \left(\frac{n^{-10}}{2(t_{\max} + \alpha)}\right)^{\frac{1}{\beta}}$$

To prove 2), consider the probability of having motion of at least $\frac{s}{\varrho}$ indices whenever the object moves. This will happen if $\mu_t \geq -0.1 \frac{s}{\varrho}$ for $t = 1, \dots, t_{\max}$. Also, if $\mu_t \leq .1 \frac{s}{\varrho}$, then the

object will move by fewer than $1.2\frac{s}{\varrho}$ indices. Using a standard Gaussian tail bound,

$$\begin{aligned}
\mathbb{P}\left(|\mu_t| \leq \frac{0.1s}{\varrho}\right)^{t_{\max}} &\geq \left(1 - \frac{2 \exp\left(\frac{-\left(\frac{0.1s}{\varrho}\right)^2}{2\sigma^2}\right)}{\frac{0.1s}{\varrho\sigma}\sqrt{2\pi}}\right)^{t_{\max}} \\
&= \left(1 - \frac{20\sigma\varrho \exp\left(\frac{-\left(\frac{0.1s}{\varrho}\right)^2}{2\varrho^2\sigma^2}\right)}{s\sqrt{2\pi}}\right)^{t_{\max}} \\
&\geq \left(1 - \frac{20\sqrt{\frac{1}{4000\log(n)}} \exp(-20\log(n))}{\sqrt{2\pi}}\right)^{t_{\max}} \\
&\geq 1 - t_{\max} \frac{20}{\sqrt{2\pi}} \sqrt{\frac{1}{4000\log(n)}} \exp(-20\log(n)) \\
&= 1 - t_{\max} \frac{1}{\sqrt{20\pi}} \frac{n^{-20}}{\sqrt{\log(n)}} \\
&\geq 1 - \frac{1}{\sqrt{20\pi}} \frac{n^{10-20}}{\sqrt{\log(n)}} \\
&\geq 1 - \frac{n^{-10}}{2}.
\end{aligned}$$

For simplicity we assume $n \geq 2$, so that $\sqrt{20\pi\log(n)} \geq 2$.

Finally, by the union bound

$$\mathbb{P}\left((2) \text{ and } 3) \text{ hold}\right) = 1 - \mathbb{P}\left((2) \text{ or } 3) \text{ does not hold}\right) \geq 1 - 2\frac{n^{-10}}{2} = 1 - n^{-10}$$

□

Proof of Lemma 3.3.4. Consider an interval of length α . The assumption that the maximum motion each time is by at at most $\varrho_2 s$ indices with $\varrho_2 s \alpha \leq n$ ensures that once an index is removed from the support, it does not return for the next α time instants. Thus without loss of generality, we can renumber the indices so that the object starts at index 1 at the beginning of the interval.

Notice from the model that for a given choice of $\mathcal{T}_{(i)}$'s, $h(\alpha)$ is an upper bound on $h^*(\alpha)$. Thus, as long as we can construct one set of mutually disjoint $\mathcal{T}_{(i)}$'s for which (3.6) holds and for which $h(\alpha; \{\mathcal{T}_{(i),u}\}) \leq h^+\alpha$ we will be done.

Let $\mathcal{T}^{[j]}$ for $j = 1, \dots, m$ be the distinct supports (in order) of \mathbf{x}_t for $t \in \mathcal{J}_u$. That is $\mathcal{T}_t = \mathcal{T}^{[j]}$ for some j . Now define $\mathcal{T}_{(i)} = \mathcal{T}^{[i]} \setminus \mathcal{T}^{[i+1]}$ for $i = 1, \dots, m-1$, and $\mathcal{T}_{(m)} = \mathcal{T}^{[m]}$. We will show that these $\mathcal{T}_{(i)}$ are disjoint. For a $j > i$,

$$\begin{aligned}\mathcal{T}_{(i)} \cap \mathcal{T}_{(j)} &= (\mathcal{T}^{[i]} \setminus \mathcal{T}^{[i+1]}) \cap (\mathcal{T}^{[j]} \setminus \mathcal{T}^{[j+1]}) \\ &= \mathcal{T}^{[i]} \cap (\mathcal{T}^{[i+1]})^C \cap \mathcal{T}^{[j]} \cap (\mathcal{T}^{[j+1]})^C\end{aligned}$$

The model assumes the support moves in a single direction, so, for a $j > i$, if $\mathcal{T}^{[j]}$ intersects with $\mathcal{T}^{[i]}$ then it also necessarily intersects with $\mathcal{T}^{[i+1]}$. Thus, $\mathcal{T}^{[i]} \cap (\mathcal{T}^{[i+1]})^C \cap \mathcal{T}^{[j]} = \emptyset$.

Next we show that $\mathcal{T}^{[j]} \subseteq \mathcal{T}_{(j)} \cup \dots \cup \mathcal{T}_{(j+\varrho-1)}$. We use the fact that $\mathcal{T}^{[j]} \cap \mathcal{T}^{[j+\varrho]} = \emptyset$.

$$\begin{aligned}\mathcal{T}^{[j]} &= [\mathcal{T}^{[j]} \cap (\mathcal{T}^{[j+1]})^c] \cup [\mathcal{T}^{[j]} \cap \mathcal{T}^{[j+1]} \cap (\mathcal{T}^{[j+2]})^c] \cup \dots \cup [\mathcal{T}^{[j]} \cap \dots \cap (\mathcal{T}^{[j+\varrho]})^c] \\ &= [\mathcal{T}_{(j)}] \cup [\mathcal{T}^{[j]} \cap \mathcal{T}_{(j+1)}] \cup \dots \cup [\mathcal{T}^{[j]} \cap \mathcal{T}^{[j+1]} \cap \dots \cap \mathcal{T}_{(j+\varrho-1)}] \\ &\subseteq \mathcal{T}_{(j)} \cup \mathcal{T}_{(j+1)} \cup \dots \cup \mathcal{T}_{(j+\varrho-1)}\end{aligned}$$

Notice that by construction, $\mathcal{T}^{[j]} \cap \mathcal{T}_{(j)} \neq \emptyset$. This along with the fact that the support changes at least once every β time instants and the fact that once an index is removed from the support, it does not return for the next α time instants, implies that $h(\alpha) \leq \beta$ for this choice of $\mathcal{T}_{(i)}$'s. Since $h^*(\alpha) \leq h(\alpha)$ and since $h^+ = \beta/\alpha$ we are done. \square

Proof of Lemma 3.3.5. For this model, we choose $\mathcal{T}_{(j)} = [(j-1)s+1, js]$ and let $\rho = 2$. The assumption that $\alpha \leq \frac{n}{m}$ ensures that no indices re-enter the support after being removed in an interval \mathcal{J}_u (which has length α). Because the support moves down by at least one index at every time t and no indices are re-visited, $h^*(\alpha) \leq s$. Therefore $h^+ = \frac{s}{\alpha}$. The assumptions then imply that $\rho^2 h^+ \leq .0024$, which satisfies Corollary 3.3.2. \square

3.C Bounding $\zeta_{k,\text{new}}^+$ [For the purposes of review]

Recall that $\zeta_{k,\text{new}}^+ := \frac{b_{\mathcal{H},k}}{b_{\mathcal{A}} - b_{\mathcal{A},\perp} - b_{\mathcal{H},k}}$.

Let $\epsilon = 0.05c_{\text{new}}\zeta\lambda^-$. Divide the numerator and denominator by λ_{new}^- and use the bounds

$f = \frac{\lambda^+}{\lambda^-}$ and $g = \frac{\lambda_{\text{new}}^+}{\lambda_{\text{new}}^-}$ to define

$$B_k := \begin{cases} \left[\rho^2 h^+ (\phi^+)^2 (\kappa_s^+)^2 (\zeta_{j,\text{new},k-1}^+ + 2\kappa_s^+ \phi^+) \right] \left(b^2 c_{\text{new}} \frac{\gamma_{\text{new}}^2}{\lambda_{\text{new}}^-} + (1-b^2)g \right) + & k=1 \\ \left[2\rho^2 h^+ (\phi^+)^2 \zeta_{j,*}^+ + \zeta_{j,*}^+ \frac{2\kappa_s^+ \phi^+}{\alpha(1-b^2)} \right] \left(b^2 \sqrt{rc_{\text{new}}} \frac{\gamma_{\text{new}}}{\lambda_{\text{new}}^-} \right) & \\ \left[\rho^2 h^+ (\phi^+)^2 \zeta_{j,\text{new},k-1}^+ + 2\sqrt{\rho^2 h^+} \phi^+ \right] \left(b^2 c_{\text{new}} \frac{\gamma_{\text{new}}^2}{\lambda_{\text{new}}^-} + (1-b^2)g \right) + & k \geq 2 \\ \left[2\rho^2 h^+ (\phi^+)^2 \zeta_{j,*}^+ + \zeta_{j,*}^+ \frac{2\sqrt{\rho^2 h^+} \phi^+}{\alpha(1-b^2)} \right] \left(b^2 \sqrt{rc_{\text{new}}} \frac{\gamma_{\text{new}}}{\lambda_{\text{new}}^-} \right) & \end{cases}$$

$$C_k := \begin{cases} \left[\rho^2 h^+ (\phi^+)^2 (\zeta_{j,*}^+) r + 2\phi^+ (\zeta_{j,*}^+) r + 2(\zeta_{j,*}^+) r \right] \left(b^2 r \frac{\gamma^2}{\lambda_{\text{new}}^-} + (1-b^2)f \right) + & k=1 \\ \left[\frac{2\phi^+ r}{\alpha(1-b^2)} + \frac{2r}{\alpha(1-b^2)} \right] \left(b^2 \sqrt{rc_{\text{new}}} \frac{\gamma_{\text{new}}}{\lambda_{\text{new}}^-} \right) + 0.25 & \\ \left[\rho^2 h^+ (\phi^+)^2 (\zeta_{j,*}^+) r + 2\sqrt{\rho^2 h^+} \phi^+ (\zeta_{j,*}^+) r + 2(\zeta_{j,*}^+) r \right] \left(b^2 r \frac{\gamma^2}{\lambda_{\text{new}}^-} + (1-b^2)f \right) + & k \geq 2 \\ \left[\frac{2\sqrt{\rho^2 h^+} \phi^+ r}{\alpha(1-b^2)} + \frac{2r}{\alpha(1-b^2)} \right] \left(b^2 \sqrt{rc_{\text{new}}} \frac{\gamma_{\text{new}}}{\lambda_{\text{new}}^-} \right) + 0.25 & \end{cases}$$

$$D_k := 1 - (\zeta_{j,*}^+)^2 - b^2 - \left[2\zeta_{j,*}^+ \frac{1}{\alpha(1-b^2)} \right] \left(b^2 \sqrt{rc_{\text{new}}} \frac{\gamma_{\text{new}}}{\lambda_{\text{new}}^-} \right) -$$

$$\left[(\zeta_{j,*}^+)^2 \right] \left(b^2 r \frac{\gamma^2}{\lambda_{\text{new}}^-} + (1-b^2)f \right) -$$

$$\zeta_{j,\text{new},k-1}^+ B_k - c_{\text{new}} \zeta C_k - 0.15$$

Then,

$$\zeta_{\text{new},k}^+ \leq \zeta_{\text{new},k-1}^+ \frac{B_k}{D_k} + c_{\text{new}} \zeta \frac{C_k}{D_k}.$$

It is not difficult to see that B_k, C_k, D_k are increasing functions of $\zeta_{\text{new},k-1}^+$ and of r, f, g , and ζ . Since $\eta \geq 1$, so $\eta f \geq f$. Thus $b^2 \eta f + (1-b^2)f$ is increasing in b and so B_k, C_k, D_k are also increasing in b .

Using the bounds assumed in Theorem 3.2.15 and since $\zeta_{\text{new},0}^+ = 1$, we can get that $\zeta_1^+ \leq 0.15$. Using this and the fact that B_k, C_k, D_k are increasing functions of $\zeta_{\text{new},k-1}^+$, by induction, we can show that $\zeta_{\text{new},k}^+ \leq \zeta_{\text{new},k-1}^+$ and thus, for all $k \geq 1$, $\zeta_{\text{new},k}^+ \leq 0.15$.

Using $\zeta_{\text{new},k}^+ \leq 0.15$ and the bounds assumed in Theorem 3.2.15, we can show that

$$\zeta_{k,\text{new}}^+ \leq 0.4\zeta_{j,\text{new},k-1}^+ + 0.5c_{\text{new}}\zeta$$

Thus,

$$\begin{aligned} \zeta_{k,\text{new}}^+ &\leq 0.4\zeta_{j,\text{new},k-1}^+ + 0.5c_{\text{new}}\zeta = \zeta_0^+(0.4)^k + \sum_{i=0}^{k-1} (0.4)^i (0.5)c_{\text{new}}\zeta \\ &\leq \zeta_0^+(0.4)^k + \sum_{i=0}^{\infty} (0.4)^i (0.5)c_{\text{new}}\zeta \\ &\leq 0.4^k + 0.84c_{\text{new}}\zeta \end{aligned}$$

References

- [1] J. Wright and Y. Ma, “Dense error correction via l1-minimization,” *IEEE Trans. on Info. Th.*, vol. 56, no. 7, pp. 3540–3560, 2010.
- [2] E. J. Candès, X. Li, Y. Ma, and J. Wright, “Robust principal component analysis?” *Journal of ACM*, vol. 58, no. 3, 2011.
- [3] V. Chandrasekaran, S. Sanghavi, P. A. Parrilo, and A. S. Willsky, “Rank-sparsity incoherence for matrix decomposition,” *SIAM Journal on Optimization*, vol. 21, 2011.
- [4] Y. Chen, A. Jalali, S. Sanghavi, and C. Caramanis, “Low-rank matrix recovery from errors and erasures,” in *Information Theory Proceedings (ISIT), 2011 IEEE International Symposium on*. IEEE, 2011, pp. 2313–2317.
- [5] H. Xu, C. Caramanis, and S. Sanghavi, “Robust pca via outlier pursuit,” *IEEE Tran. on Information Theory*, vol. 58, no. 5, 2012.
- [6] M. B. McCoy and J. A. Tropp, “Sharp recovery bounds for convex demixing, with applications,” *arXiv:1205.1580*.
- [7] V. Chandrasekaran, B. Recht, P. A. Parrilo, and A. S. Willsky, “The convex geometry of linear inverse problems,” *Foundations of Computational Mathematics*, no. 6, 2012.
- [8] M. Tao and X. Yuan, “Recovering low-rank and sparse components of matrices from incomplete and noisy observations,” *SIAM Journal on Optimization*, vol. 21, no. 1, pp. 57–81, 2011.
- [9] F. D. L. Torre and M. J. Black, “A framework for robust subspace learning,” *International Journal of Computer Vision*, vol. 54, pp. 117–142, 2003.

- [10] M. Mardani, G. Mateos, and G. B. Giannakis, "Dynamic anomalography: Tracking network anomalies via sparsity and low rank," *Selected Topics in Signal Processing, IEEE Journal of*, vol. 7, no. 1, pp. 50–66, 2013.
- [11] J. A. Tropp, "User-friendly tail bounds for sums of random matrices," *Foundations of Computational Mathematics*, vol. 12, no. 4, 2012.
- [12] C. Qiu, N. Vaswani, B. Lois, and L. Hogben, "Recursive robust pca or recursive sparse recovery in large but structured noise," *IEEE Trans. Info. Th.*, Aug. 2014, shorter versions in ICASSP 2013 and ISIT 2013.
- [13] B. Lois and N. Vaswani, "A correctness result for online robust pca," *arXiv:1409.3959*.
- [14] H. Guo, C. Qiu, and N. Vaswani, "An online algorithm for separating sparse and low-dimensional signal sequences from their sum," *IEEE Trans. Sig. Proc.*, pp. 4284–4297, Aug. 2014.
- [15] J. He, L. Balzano, and J. Lui, "Online robust subspace tracking from partial information," *arXiv:1109.3827 [cs.IT]*, 2011.
- [16] G. Mateos and G. B. Giannakis, "Robust pca as bilinear decomposition with outlier-sparsity regularization," *Signal Processing, IEEE Transactions on*, vol. 60, no. 10, pp. 5176–5190, 2012.
- [17] D. Hsu, S. Kakade, and T. Zhang, "Robust matrix decomposition with sparse corruptions," *IEEE Trans. Info. Th.*, Nov. 2011.
- [18] J. Feng, H. Xu, and S. Yan, "Online robust pca via stochastic optimization," in *Advances in Neural Information Processing Systems 26*, C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Weinberger, Eds. Curran Associates, Inc., 2013, pp. 404–412. [Online]. Available: <http://papers.nips.cc/paper/5131-online-robust-pca-via-stochastic-optimization.pdf>
- [19] J. Feng, H. Xu, S. Mannor, and S. Yan, "Online pca for contaminated data," in *Advances in Neural Information Processing Systems 26*, C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Weinberger, Eds. Curran Associates, Inc., 2013, pp. 764–772. [Online]. Available: <http://papers.nips.cc/paper/5135-online-pca-for-contaminated-data.pdf>
- [20] J. Zhan and N. Vaswani, "Robust pca with partial subspace knowledge," *arXiv:1403.1591 [cs.IT]*, 2014.
- [21] S. Li and H. Qi, "Recursive low-rank and sparse recovery of surveillance video using compressed sensing," in *Proceedings of the International Conference on Distributed Smart Cameras*, ser. ICDSC '14. New York, NY, USA: ACM, 2014, pp. 1:1–1:6. [Online]. Available: <http://doi.acm.org/10.1145/2659021.2659029>
- [22] B. Nadler, "Finite sample approximation results for principal component analysis: A matrix perturbation approach," *The Annals of Statistics*, vol. 36, no. 6, 2008.
- [23] E. Candes, "The restricted isometry property and its implications for compressed sensing," *Compte Rendus de l'Academie des Sciences, Paris, Serie I*, pp. 589–592, 2008.

- [24] M. Muselli, “On convergence properties of pocket algorithm,” *Neural Networks, IEEE Transactions on*, vol. 8, no. 3, pp. 623–629, May 1997.
- [25] D. Hsu, S. M. Kakade, and T. Zhang, “Robust matrix decomposition with sparse corruptions,” *Information Theory, IEEE Transactions on*, vol. 57, no. 11, pp. 7221–7234, 2011.
- [26] C. Davis and W. M. Kahan, “The rotation of eigenvectors by a perturbation. iii,” *SIAM Journal on Numerical Analysis*, Mar. 1970.
- [27] R. Horn and C. Johnson, *Matrix Analysis*. Cambridge Univ. Press, 1985.

CHAPTER 4. GENERAL CONCLUSIONS

In this work correctness results for two versions the ReProCS algorithm were proved. In Chapter 2, results were obtained for both the online matrix completion problem and the online robust PCA problem. It was shown that if 1) the algorithm parameters were set appropriately, 2) the ℓ_t are independent over t and lie in a dense low-dimensional subspace that changes slowly over time, and 3) the support of \mathbf{x}_t changes enough with t , then the ReProCS algorithm will accurately recover \mathbf{x}_t and ℓ_t at each time t .

Chapter 3 analyzed a modification of the ReProCS algorithm that included a cluster PCA step. This allowed for directions to be removed from the estimate of the subspace where the ℓ_t lie. Also in this paper, the assumption that the ℓ_t be independent was relaxed, and an autoregressive model was assumed on the coefficients \mathbf{a}_t .

In ongoing work, we have been able to remove the assumption that the eigenvalue clusters are known a priori. Using the same proof techniques (Weyl's inequality and the Matrix Hoeffding inequality), we can show that for each i , $\lambda_i \left(\sum_t \hat{\ell}_t \hat{\ell}_t' \right)$ is close to $\lambda_i \left(\sum_t \ell_t \ell_t' \right)$. Thus, the clusters can be detected automatically.

As mentioned in Chapter 2, one direction for new work would be to study the under-sampled case $\mathbf{m}_t = \mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \ell_t$ where \mathbf{A}_t and \mathbf{B}_t are matrices with more columns than rows. A partial result for the case where \mathbf{B}_t does not change with t was proved in [1]. Using the new techniques introduced in this thesis, it should be straightforward to extend this result to a complete correctness result. Another open problem is the noisy case: $\mathbf{m}_t = \mathbf{x}_t + \ell_t + \mathbf{w}_t$ where \mathbf{w}_t is small and bounded noise.

References

- [1] B. Lois, N. Vaswani, and C. Qiu, “Performance guarantees for undersampled recursive sparse recovery in large but structured noise,” in *GlobalSIP*, 2013.